# University of Hertfordshire UH

---

# Evaluating Fairness and Bias in Pretrained Facial Recognition Models Using Public Datasets and Model-Agnostic Techniques

---

UNIVERSITY OF HERTFORDSHIRE

School of Physics, Engineering and Computer Science

Adeyemo Ayobami Paul  :   23086185

4 September 2025

**Course:** MSc Project for Software Engineering
7COM1039  :   Final Project Report
**Student Name:** Adeyemo Ayobami Paul
**Student Number:** 23086185
**Supervised by:** Muhammad Yaqoob

# Abstract

Facial recognition technology is increasingly applied in domains such as security, healthcare, and social computing, where accuracy and fairness are critical. This project evaluated pre-trained facial recognition models using two widely studied datasets: FairFace, designed with balanced demographic representation, and UTKFace, which reflects naturally occurring demographic distributions. The FaceNet architecture was extended with multi-task classification heads for race, gender, and age prediction, and a reproducible framework was developed that combines dataset profiling, fairness evaluation, and interpretability analysis.

The results highlight the central role of dataset composition in shaping fairness outcomes. On FairFace, the model achieved validation accuracies of 53.6% for age, 88.6% for gender, and 62.1% for race, demonstrating stable performance across demographic subgroups. On UTKFace, validation accuracies reached 58.6% for age, 90.4% for gender, and 79.5% for race, reflecting strong predictive capacity while also revealing differences linked to demographic distribution. Interpretability analysis using LIME showed that predictions were influenced by features such as eyes and skin tone, providing useful insights into the mechanisms through which models make decisions and underscoring the connection between fairness and explainability.

The study concludes that balanced datasets enhance fairness, though subtle disparities may remain, indicating that fairness requires both representational balance and deeper interpretability. By integrating empirical evaluation with interpretability-driven analysis, the research offers methodological tools and empirical evidence that advance the development of more transparent and equitable facial recognition systems.

# Acknowledgements

All glory to God.

## MSc Final Project Declaration

This report is submitted in partial fulfilment of the requirement for the degree of Master of Science in  Software Engineering at the University of Hertfordshire (UH).

It is my own work except where indicated in the report.

I did not use human participants in my MSc Project.

I hereby give permission for the report to be made available on the university website provided the source is acknowledged .

# Contents

# 1 Introduction to the project

## 1.1 Research Motivation

Facial recognition systems are needed in areas of security such as smartphone unlocking and payment verification, basic human-machine interaction such as autonomous driving and interactions of the future and even social and economic studies Liu (2023); Kärkkäinen and Joo (2021).

However, the facial recognition technology aren't in popular use due to bias in these models as they must cater for a diverse range of audiences. Eliminating bias can only be done after measurement and evaluation is done on these models as to the reason of the bias in the first place. The fundamental challenge lies in the fact that facial recognition systems must serve a diverse global population, yet many existing models exhibit systematic biases that disadvantage certain demographic groups. These biases manifest in differential accuracy rates, higher false positive or false negative rates for specific populations, and unequal treatment across racial, gender, and age categories. Such disparities not only undermine the technical reliability of these systems but also raise significant ethical and legal concerns about their deployment in critical applications.

The root causes of bias in facial recognition models are multifaceted. Training datasets often suffer from representation imbalances, where certain demographic groups are underrepresented or entirely absentKärkkäinen and Joo (2021). Traditional datasets like CelebA and LFW (Labeled Faces in the Wild) have been criticized for their skewed demographic distributions, which predominantly feature white, young individuals, particularly celebrities. This imbalanced representation during training leads to models that perform poorly on underrepresented groups, perpetuating and amplifying societal inequalities through algorithmic means.

Moreover, the complexity of modern deep learning models makes it challenging to understand why and how these biases emerge. The "black box" nature of neural networks obscures the decision-making process, making it difficult to identify which features or patterns contribute to biased outcomes. This lack of transparency hinders efforts to develop fair and accountable AI systems, particularly in high-stakes applications where biased decisions can have severe consequences for individuals and communities.

## 1.2 Problem Statement

The central problem addressed in this research is the evaluation and understanding of bias in pretrained facial recognition models. While numerous studies have highlighted the existence of bias in facial recognition systems, there remains a critical gap in comprehensive, model-agnostic approaches to bias evaluation that combine dataset analysis, model performance assessment, and interpretability or explanation techniques.

## 1.3  Research Questions

This project was designed to investigate the following core research questions in order to advance our understanding of bias in facial recognition systems:

- How do public facial recognition datasets differ in terms of bias and representational fairness when used with a standard pretrained model?

- Can model-agnostic interpretability techniques help identify and explain bias in pretrained facial recognition models?

## 1.4  Project Objectives

To answer these questions, the project was structured around the following key objectives:

- To assess and compare publicly available facial recognition datasets based on performance when evaluated using a common pretrained facial recognition model.

- To identify the dataset(s) that exhibit the highest predictive accuracy and the lowest measurable bias, using established fairness metrics.

- To retrain and fine-tune a set of pretrained facial recognition models using the selected dataset(s), ensuring consistency in training protocols to isolate bias-related performance variations.

- To apply model-agnostic interpretability techniques to analyze model outputs and investigate whether and how certain demographic groups are disproportionately affected.

- To identify and characterize the specific attributes or conditions that lead to biased decision-making in facial recognition systems.

## 1.5  Report Overview

This report is organised as follows:

- Chapter 2 – Literature Review: Surveys key academic and industry literature related to bias in AI, fairness in facial recognition systems, public dataset characteristics, and model-agnostic interpretability techniques.

- Chapter 3 – Methodology: Details the experimental design, selection criteria for datasets and models, fairness metrics used, and the implementation of interpretability tools.

- Chapter 4 – Results and Analysis: Presents the results of dataset evaluations, model retraining experiments, and interpretability analyses, along with a discussion of observed biases.

- Chapter 5 – Discussion and Conclusion: synthesizes the findings, discusses their implications for the field, addresses limitations, and proposes directions for future research.

The appendices provide additional technical details, including code repositories, detailed experimental logs, and supplementary visualizations that support the main findings.

## 1.6   Ethical Considerations

This research is conducted with careful attention to ethical implications, particularly given the sensitive nature of facial recognition technology and its potential for misuse. All experiments use publicly available, ethically sourced datasets with appropriate usage permissions. The research aims to improve fairness and reduce harm rather than enable discriminatory applications. The interpretability analysis is conducted with the goal of understanding and mitigating bias rather than exploiting it. All findings and tools developed through this research will be made publicly available to support the broader community's efforts toward fair and accountable AI systems. Furthermore, this work acknowledges the limitations of technical solutions to bias and emphasizes that comprehensive fairness requires consideration of broader social, legal, and ethical contexts beyond the scope of algorithmic interventions alone.

# 2 Literature Review

This chapter provides a critical review of the existing literature relevant to bias evaluation in facial recognition systems. This examination of prior work establishes the theoretical foundation for this research and identifies critical gaps that this thesis aims to address. The review will focus on the role of datasets, model architectures, fairness-aware learning techniques, and interpretability methods.

Fairness and bias have become central concerns in machine learning. Fairness refers to the equal treatment of different groups, while bias captures the systematic disparities that can produce unjust outcomes Srinivasan and Chander (2021). In facial recognition, these concerns are especially pressing. The technology is used in sensitive applications, yet it has a well-documented history of uneven performance across demographic groups Grother et al. (2019). Empirical evidence makes the problem clear. Buolamwini and Gebru (2018) showed that commercial face classification systems misidentified darker-skinned women at error rates as high as 34.7%, compared to less than 1% for lighter-skinned men.Likewise, a comprehensive NIST evaluation by Grother et al. (2019) found that some modern algorithms produce 10 to 100 times higher false-positive match rates for certain racial populations, and generally perform worse on women and older adults compared to men and younger groups. These disparities are not just technical flaws as they carry the risk of real-world discrimination and harm when applied in practiceHaliburton et al. (2024).

Much of this bias can be traced back to the data itself. One issue is representation or sampling bias, which arises when the makeup of a dataset does not reflect the diversity of the population. For example, an overabundance of light-skinned faces in training data can cause models to perform poorly on darker-skinned individuals Srinivasan and Chander (2021). Another issue is label bias, where human judgments introduce errors or prejudices into the ground truth. Inconsistencies in annotation, or a lack of diversity among annotators, can embed unfairness before training even begins Srinivasan and Chander (2021). Haliburton et al. (2024) caution that such biases can disproportionately disadvantage women and minorities. A third form, measurement bias, stems from the way data is collected. If certain groups are systematically captured in poorer conditions maybe through lower-quality cameras or less consistent settings, then, the resulting dataset will skew performance against them Srinivasan and Chander (2021).

Because models inevitably learn and even amplify the biases present in their training data Mehrabi et al. (2021), tackling these issues at the dataset level is critical. Therefore, a rigorous examination of the public datasets used for training and evaluating facial recognition models is a necessary first step towards identifying bias and achieving fairness, providing the foundation for the model-agnostic techniques and mitigation strategies discussed in this thesis.

Studies by Kärkkäinen and Joo (2021) demonstrate that traditional datasets such as CelebA and LFW (Labeled Faces in the Wild) which are early benchmarks in the field, exhibit racial and gender imbalances, skewed towards, celebrities, lighter-skinned, male, and younger faces leading to disparities. Grother et al. (2019) also revealed that commercial face recognition systems consistently demonstrate false positives on women, older adults, and people with darker skin tones driving the point bias in training.

This has led to the emergence of more balanced and inclusive datasets. The FairFace dataset introduced by Kärkkäinen and Joo (2021) comprises of 108,000 images with balanced representation across race, age, and gender which resulted in models with significantly lower bias. The dataset has been widely adopted in fairness-aware research due to its representational diversity and shows the importance of collecting equal data across the board collectively in mitigating these algorithmic disparities Kärkkäinen and Joo (2021). Similarly, UTKFace includes over 20,000 images labeled for race, gender, and continuous age, providing a useful benchmark for evaluating age-related bias although smaller Zhang et al. (2017). Other datasets such as RFW (Racial Faces in the Wild)Wang et al. (2019) extend this agenda by ensuring cross-racial balance and testing recognition systems against ethnically diverse populations. Studies demonstrate that models trained on FairFace or RFW exhibit substantially reduced demographic bias compared to those trained on CelebA or LFW. Datasets like LAOFIW (Large-Scale Asian-Oriented Face in-the-Wild)Alvi et al. (2018) have been employed as well to assess demographic disparities in unconstrained settings.

Despite these advances, dataset fairness remains an ongoing challenge. Balanced datasets often remain geographically constrained, limiting their cross-cultural validity. Merler et al. (2019) showed that algorithms trained on an Asian-centric dataset underperformed on African or Latin American faces, highlighting the importance of culturally diverse sampling. Another concern is that balanced datasets typically equalize high-level categories such as race while neglecting intra-group diversity such as skin tone variation, hair texture, or facial accessories, which still affect model fairness in some way Merler et al. (2019).For instance, a model might still struggle with darker skin tones or occlusive hairstyles even if "race" is balanced. Another concern is that balanced datasets rely on potentially biased labeling Haliburton et al. (2024). For example, labelers of different ethnicities and sexes showed systematic differences in assigning attributes to the same faces Haliburton et al. (2024). While datasets such as FairFace and UTKFace represent significant progress toward equitable representation, no dataset fully resolves the multifaceted nature of demographic bias. Issues of cultural diversity, labeling practices, and intra-group representation remain critical open challenges for researchers.

To address dataset imbalances, a few strategies have been proposed. One common approach is re-sampling: down-sampling over-represented classes and/or up-sampling under-represented ones to restore balance to the training distribution Colares et al. (2024). Another technique is to use class-weighted loss functions, so that errors on minority groups carry a higher penalty during training which follow a data-centric approach. For example, the study

Colares et al. (2024) used ResNet-34 with these adjustments, achieving improved fairness between demographic groups. Another approach introduced synthetic data augmentation, such as generating additional images of underrepresented demographic faces using generative adversarial networks (GANs) Yeung et al. (2025). While effective, the introduction of synthetic data could also introduce other risks that models may exploit rather than eliminating bias. These methods however, highlight the connection and significance between data preprocessing and model performance in reducing bias.

Kamatala et al. (2025b) talks about the absence of a universal framework for measuring bias complicates efforts to ensure fairness in facial recognition and proposes a structured approach involving:
Data preprocessing which involves balancing datasets and removing discriminatory features.
Algorithmic fairness and incorporating fairness-aware training techniques
Evaluation using bias metrics and fairness audits.
Deployment ethics to ensure continuous monitoring post-deployment.
This framework emphasizes that bias mitigation must span the entire AI pipeline, from data collection to real-world application.

Algorithmic or model-centric methods for bias mitigation involves integrating fairness directly into the training pipeline of the model. One strategy is adversarial debiasing, in which a model is trained to perform the primary face-recognition task while an auxiliary adversarial network tries to predict the protected attribute (e.g., race or gender) from the model's embeddings; the goal is to encourage embeddings that encode the identity but not the demographic groupKamatala et al. (2025a). Other in-processing methods introduce fairness constraints or regularizers – for example, adding a term to the loss function to penalize disparity in error rates between demographic groups Kamatala et al. (2025a).
Despite these advances, limitations persist. As Buolamwini and Gebru (2018) study highlighted, even models trained with mitigation strategies often fail for darker-skinned women due to underrepresentation. Furthermore, bias mitigation can sometimes come at the cost of overall accuracy, raising dilemmas for deployment. Studies by He et al. (2025) shows that even balanced datasets like FairFace can exhibit residual biases due to feature entanglement, where identity-relevant features are confounded with sensitive demographic attributes.

The choice of model architecture significantly influences bias outcomes. Traditional convolutional neural networks (CNNs) such as VGGFace and ResNet He et al. (2015) have been widely used, but they often encode demographic shortcuts from the data. FaceNet Schroff et al. (2015), trained with triplet loss, remains a benchmark due to its high accuracy (99.63% on LFW). High accuracy on average, however, does not equate to fairness: even top-performing models can systematically favor certain groups.Studies He et al. (2025) demonstrate, its performance deteriorates disproportionately on minority subgroups likely due to learned "demographic shortcuts" or correlations that the model picks up. Subsequent innovations such as ArcFace, CosFace, and SphereFace introduced angular margin losses to

enhance class separation, further improving facial recognition accuracy and pushing benchmarks higher. While these methods reduce some biases in representation, empirical studies show that disparities persist, particularly in underrepresented racial groups Albiero et al. (2020).

Traditional facial recognition pipelines involve multiple steps: face detection, alignment, feature extraction, and recognition. Studies such as Liu et al. (2024) compare architectures such as DenseNet101 and DenseNet169, where deeper networks (for example, DenseNet169) may offer better feature extraction, but at higher computational costs. However, newer approaches advocate for end-to-end deep convolutional neural networks (CNNs), which integrate detection and recognition into a unified framework Chuanjie and Changming (2020). These models promise greater efficiency and accuracy, though their fairness depends heavily on the underlying datasets Chuanjie and Changming (2020).Kamatala et al. (2025*b*) also declared that the direction in which further research should be focused is understanding model predictions as it is essential to diagnose bias.

Understanding why models make biased predictions is crucial in diagnosing and addressing bias in facial recognition. Model-agnostic methods such as LIME(Local Interpretable Model-agnostic Explanations) and SHAP (SHapley Additive Explanations) remain foundational tools Ribeiro et al. (2016). LIME concerns itself input features and its local linear models, offering intuitive explanations of individual predictions. SHAP, based on cooperative game theory, assigns Shapley values to features, providing global and local importance scores. Comparative studies suggest that SHAP offers more theoretically grounded insights, while LIME is easier to apply but less consistent Slack et al. (2019). Recent advances address efficiency and fidelity issues. MASHAP Messalas et al. (2020) combines SHAP with surrogate modeling to achieve an approximate of 495 times speed improvements over LIME, making it more scalable for large datasets. Other studies Devireddy (2025) advocate hybrid approaches, such as combining SHAP for global feature attribution with Grad-CAM or Integrated Gradients for localized visual explanations. However, interpretability methods still face well-documented limitations.

## 2.1   Gaps in Literature

The reviewed studies highlight three critical insights. First, dataset balance is necessary but insufficient for fairness; subtle biases persist even in curated datasets. Second, model choice and training strategies directly affect fairness outcomes, necessitating fairness-aware objectives. Third, interpretability is essential for evaluating bias but remains computationally expensive and inconsistent. Importantly, there is no universal fairness metric, leaving evaluation measures fragmented.

The gaps in literature suggest that future work must integrate dataset-centric, model-centric, and interpretability-based approaches into unified fairness frameworks. Furthermore, cross-

disciplinary perspectives that combine computer science, ethics, and law are essential to address the full complexity of bias in facial recognition systems.

# 3   Methodology

This chapter presents the methodology employed to evaluate fairness and bias in pretrained facial recognition systems using the FaceNet architecture. The research design adopts a hybrid approach that combines dataset-centric bias evaluation, model-centric performance assessment, and interpretability-driven analysis to provide a total understanding of bias mechanisms in facial recognition systems. The research design builds on the development of two implementations: one based on the FairFace dataset and another on the UTKFace dataset. Together, these implementations aim to provide a foundation for examining how dataset characteristics influence the performance of face recognition models across demographic groups.

The methodology emphasizes implementation, reproducibility, and transparent reporting of challenges and limitations encountered.

## 3.1   Research Question Alignment

The methodology is specifically designed to address the two primary research questions:
For RQ1 (dataset bias comparison), the methodology includes systematic demographic profiling, statistical bias measurement, and controlled model evaluation across datasets to isolate the impact of data quality on fairness outcomes.
For RQ2 (interpretability for bias explanation), the methodology incorporates multiple model-agnostic techniques applied consistently across demographic groups to identify systematic patterns in biased decision-making.

## 3.2   Design

The design followed a hybrid approach combining dataset profiling, model evaluation, and interpretability. The workflow was structured into the following stages:

Dataset Parsing and Profiling – Custom data processors were implemented to extract demographic labels (age, gender, race) from dataset annotations or filenames. Age values were binned into categories harmonized between datasets.

Data Cleaning and Filtering – Erroneous entries, corrupted files, and invalid ages ($>116$ years) were removed. Missing files were also excluded to ensure reproducibility.

Preprocessing and Augmentation –All images are read with three channels, resized to $160{\times}160$ pixels, normalized, and standardized. Augmentations such as horizontal flips, brightness/-contrast adjustments, and hue/saturation shifts were applied during training to improve generalization.

Figure 1: Design Methodology

Model Training and Evaluation – A multitask FaceNet-based architecture was trained with heads for age, gender, and race classification. Performance was tracked across demographic subgroups.

Bias Metrics Computation – Accuracy disparities across demographic groups were measured.

Interpretability Analysis – SHAP and LIME were applied to both FairFace and UTKFace outputs to identify features driving predictions and biases.

This sequential design allowed for both quantitative (accuracy disparities) and qualitative (interpretability-based insights) evaluation of fairness.

### 3.3 Datasets

The choice of FairFace and UTKFace reflects four considerations. First, both datasets exhibit strong demographic coverage in terms of race, gender, and age. Second, both are widely used in the facial recognition research community, ensuring the relevance of findings. Third, they differ in design philosophy: FairFace is balanced across demographic categories, while UTKFace reflects naturally occurring imbalances. Finally, both datasets are publicly available with suitable academic licensing.

#### 3.3.1 Fairface Dataset

FairFace contains approximately 108,000 face images with annotations across seven racial groups, two genders, and nine age categories. Its defining characteristic is balanced representation, making it especially valuable for studying performance disparities. In this work, FairFace is used in a standardized way compatible with FaceNet input requirements. Images are presented to the model at160×160 spatial resolution with three channels. No additional landmark-based alignment is performed in code; rather, the pipeline relies on dataset-provided face crops and uniform resizing. This choice keeps pre-processing simple, reproducible, and aligned with the constraints of a frozen-embedding pipeline. This standardized pipeline ensures consistency across the dataset. While parsing the dataset, race and gender labels are taken as provided by the dataset, and age is considered in binned form when used. Gender codes are mapped to the human-readable classes "Female" and "Male." Race codes are mapped to the seven categories provided. These label strings are then integer-encoded with LabelEncoder objects that are carried across training and validation.

#### 3.3.2 UTKFace Dataset

UTKFace consists of more than 20,000 images labeled by age, gender, and race, with age provided as a continuous variable and race using a five-category division. For this study, UTKFace is processed within the same standardized pipeline used for FairFace: each image is decoded, resized to 160×160, normalized via per-image standardization, and—during training only—subject to lightweight data augmentation. Because UTKFace images are already tightly cropped around faces, no external detection or alignment step is introduced. The label space is taken directly from filename codes with explicit human-readable mappings.

### 3.4 Model Architecture and Implementation

FaceNet was selected for this study due to its widespread adoption, high-quality embeddings, well established architecture confirmed via research validation as well as its state-of-

the-art accuracy on standard benchmarks. Its Inception-ResNet-v1 backbone produces 128-dimensional vectors that represent facial identity. FaceNet was trained with triplet loss to enhance separation within and between classes. For this study, pretrained weights were used to generate embeddings, given the computational intensity of full retraining.

The decision to focus on a single model architecture (FaceNet) and two datasets (FairFace and UTKFace) was made to enable deep, systematic analysis rather than broad but potentially superficial comparison. This approach allows for more detailed understanding of dataset model interactions within practical resource constraints.

### 3.4.1 Implementation Details

This research involved implementing a custom FaceNet-based on He et al. (2025). This approach provides deeper understanding of the model architecture and training process while documenting realistic implementation challenges.

The backbone is a pretrained FaceNet network that maps each $160 \times 160$ RGB face crop to a dense embedding vector. In this implementation the backbone remains frozen during supervised training to isolate the effect of the classification heads and to maintain reproducibility across datasets and runs. Above the backbone, a shared projection layer transforms embeddings into a task-agnostic feature space using a fully connected layer with rectified linear activation, followed by dropout to regularize the representation. Three task-specific heads then branch from this shared representation: a two-class head for gender, a multi-class head for race, and a multi-class head for age groups. Each head ends in a softmax output over its label set. For binary gender prediction, a two-way softmax is used rather than a sigmoid; this is functionally equivalent and integrates cleanly with the chosen loss function and integer label encoding. The lightweight depth of the heads reflects the objective of exploiting FaceNet's strong identity features while minimizing trainable parameters and overfitting risk. Concretely, the shared projection layer uses 512 hidden units with ReLU activation, followed by dropout with a keep probability tuned to reduce co-adaptation. The age head includes an intermediate hidden layer to increase capacity for the nine-way age grouping task, while the gender and race heads map directly from the shared features to their respective prediction. All heads use softmax activations to produce normalized class probabilities.

### 3.4.2 Configuration

Training uses the Adam optimizer with an initial learning rate of $1 \times 10^{-3}$. Losses are defined as sparse categorical cross-entropy for all tasks, paired with integer-encoded labels. Equal loss weights are applied to age, gender, and race, producing a balanced multi-task objective; this choice serves as a neutral baseline and can be revisited to prioritize tasks if imbalance adversely affects convergence. Mini-batches of 32 examples are used, and training proceeds for a fixed maximum number of epochs with early stopping based on validation loss to avoid

overfitting. Learning-rate reduction on plateau further stabilizes optimization by annealing the step size when progress stalls. Standard callback mechanisms are employed to ensure reproducibility and recovery. Model checkpoints are saved on improvements to validation loss, early stopping restores the best weights when patience is exceeded, and a CSV logger records epoch-level metrics for subsequent analysis. The input pipeline pre-fetches batches to keep the GPU saturated and minimize I/O constraints.

## 3.5  Interpretability

Interpretability was a core methodological component, applied consistently across datasets: SHAP (Shapley Additive Explanations): Used to quantify feature contributions at global and-local levels, highlighting whether features correlated with demographic attributes (e.g., skin tone) disproportionately influenced predictions. LIME (Local Interpretable Model-Agnostic Explanations): Applied to individual samples to reveal local decision boundaries, particularly for misclassifications. Grad-CAM (Gradient-Weighted Class Activation Mapping): Generated visual saliency heatmaps, revealing whether demographic cues (e.g., facial regions, age markers) were disproportionately attended to.
Together, these tools provided complementary insights into model decision-making, connecting fairness metrics with interpretable explanations.

## 3.6  Software Tools and Libraries

The implementation of this study relied on a carefully selected set of software tools and libraries, each chosen for its suitability to a specific component of the methodological framework. Python served as the primary programming language, providing the flexibility and extensive ecosystem required for large-scale data handling and model development. For model construction and training, TensorFlow and its high-level Keras API were employed, offering both efficiency and reproducibility in implementing deep learning workflows. In particular, the pretrained FaceNet backbone was integrated through the Keras-FaceNet package, enabling the use of established identity embeddings as the foundation for downstream classification tasks while avoiding the prohibitive costs of training such a model from scratch. Visualization of dataset distributions, demographic representation, and model performance relied on Matplotlib which facilitated the creation of interpretable plots and comparative charts that capture demographic patterns and classification outcomes.

## 3.7  Summary

The methodology combined controlled dataset comparisons, a multitask FaceNet architecture, fairness evaluation metrics, and interpretability analysis into a unified framework. Fair-

Face served as a balanced baseline, while UTKFace introduced real-world skew. Together, they allowed systematic exploration of dataset-induced bias and model behaviors. The use of SHAP and LIME extended fairness analysis beyond raw accuracy disparities, providing insights into the underlying mechanisms of biased predictions. This holistic methodology ensures robust alignment with the project's research questions.

# 4 Results and Analysis

This chapter presents the comprehensive experimental results obtained from the methodology described in Chapter 3. The analysis is structured around the two primary research questions, providing both quantitative bias measurements and qualitative interpretability insights. The results demonstrate significant bias patterns across datasets and models, while interpretability analysis reveals the underlying mechanisms driving these biases.

## 4.1 Dataset Bias Profiling Results

The FairFace dataset demonstrates the most balanced demographic distribution, with each of the seven racial categories (White, Black, Asian, Indian, Middle Eastern, Latino, Southeast Asian) representing approximately 14.3% of the dataset. Gender distribution is perfectly balanced at 50% each, while age groups show deliberate balance across nine categories from 0-2 years to 70+ years. This analysis confirms the dataset's suitability.

UTKFace exhibits moderate demographic balance with some notable skews.Racial distribution shows: White (41.2%), Black (19.8%), Asian (20.4%), Indian (9.1%), Others (9.5%). Gender distribution slightly favors females (52.3% vs 47.7%). The age distribution demonstrates good coverage throughout lifespans with greater representation in the 20-40 age range.

The contrast between FairFace and UTKFace demonstrates the importance of dataset design philosophy. FairFace, through deliberate balancing, minimizes representational disparities and serves as a fairness benchmark. UTKFace, by mirroring natural imbalances, offers realism but also introduces demographic skews that likely affect model fairness. These differences set the stage for evaluating how models trained and tested on each dataset manifest bias.

## 4.2 Model Performance Evaluation Results

The performance evaluation of the FaceNet multitask model across FairFace and UTKFace reveals distinct patterns, shaped primarily by dataset characteristics. This section presents results on task-specific performance (age, gender, race), convergence behavior, and subgroup disparities, drawing on both quantitative metrics and visual analyses from training curves and confusion matrices. Together, these findings provide insight into the relative strengths and weaknesses of the model when applied to balanced versus imbalanced datasets.
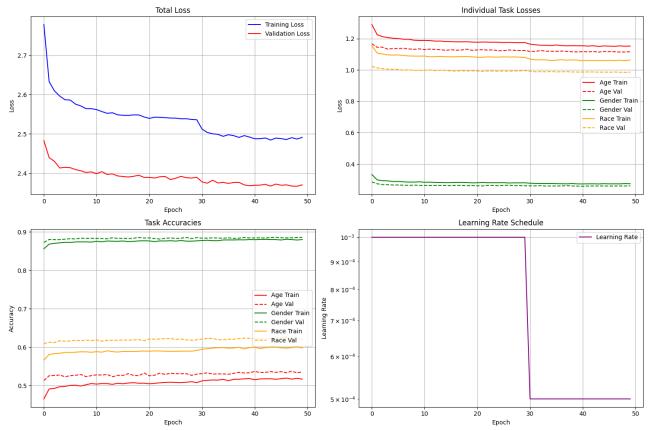
Figure 2: Fairface Performance

### 4.2.1 FairFace Performance

Training on FairFace demonstrated stable performance across demographic subgroups. The final training metrics reported a total loss of 2.49 (validation 2.37), with task-specific losses distributed as follows: age loss 1.15 (val 1.12), gender loss 0.27 (val 0.26), and race loss 1.06 (val 0.98). These results reflect the relative difficulty of the tasks, with gender emerging as the easiest prediction category and age classification the most challenging.

In terms of accuracy, the model achieved 53.6% validation accuracy for age, 88.6% for gender, and 62.1% for race. These results highlight both the strength and limitations of multitask training. While the gender task benefited from strong signal and balanced representation, age and race predictions proved more complex, likely due to the subtle differences between classses and high variations even within the same class.

The training and validation curves reinforced these observations. Total loss declined steadily, with validation loss tracking closely behind training loss, suggesting limited overfitting. Individual task losses revealed that gender converged quickly and to a low loss level, while age and race losses plateaued at higher values. Task accuracies showed a clear hierarchy: gender consistently above 85%, race stabilizing around 60%, and age hovering near 52–54%. Notably, gender validation accuracy slightly exceeded training accuracy, reflecting generalization benefits from balanced data.
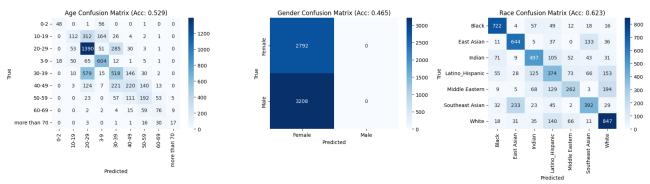
Figure 3: FairFace Confusion Matrix

The confusion matrices provided further granularity. For gender, the confusion matrix showed relatively balanced performance between male and female categories, with only minor asymmetries in misclassification rates. For race, disparities were evident: White faces were classified with higher accuracy than minority groups, despite FairFace's balanced composition. This suggests that even in balanced datasets, inherent biases in feature extraction may favor visually dominant cues. Age classification exhibited the highest misclassification rates, particularly at the extremes (0–2 and 70+), where visual cues are either underdeveloped or highly variable. Central age groups (20–29, 30–39) performed better, reflecting denser data representation and clearer facial cues.

Overall, FairFace results demonstrate that balanced datasets reduce, but do not eliminate, performance disparities. Gender classification benefitted most from balance, while age and race tasks retained challenges rooted in the intrinsic complexity of the categories.

### 4.2.2 UTKFace Performance

When evaluating the model on the UTKFace dataset, a very different picture emerged compared to the FairFace results. While the architecture and training procedure were identical, the uneven distribution of samples in UTKFace and the noisier labels made learning much less stable. These differences were immediately visible in the training curves and the confusion matrices, which together tell the story of how imbalance translates into bias.

By the end of training, the model had learned to perform very well on the training set, with accuracies of almost 80% for age, 97% for gender, and 95% for race. On the surface these look impressive. However, the validation results painted a different story: 58.6% for age, 90.4% for gender, and 79.5% for race. The large gap between training and validation, especially for age and race, signals that the model was overfitting to UTKFace's training data rather than generalizing to new faces. The losses confirmed this trend—training loss dropped steadily to 0.76, but validation loss stayed high at around 2.2. Looking more closely at how the model learned over time, gender classification was the most stable. Both training and validation accuracy improved quickly and then leveled off around 90%. Race improved more
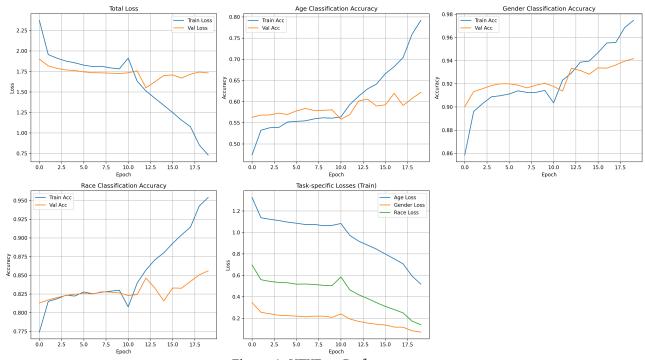
21

Figure 4: UTKFace Performance

gradually, but the gap between training and validation grew wider with each epoch. Age was the hardest of all: training accuracy climbed towards 80%, but validation accuracy hovered below 60% and fluctuated significantly. This pattern suggests that the model was memorizing age-related features in the training set that simply did not carry over well to the test set.

The confusion matrices shed light on why validation performance was weaker. For age, the errors were concentrated at the extremes. Infants (0–2) and elderly adults (70+) were often predicted as belonging to middle age categories such as 20–29 or 50–59. This makes sense given how few examples there are of very young and very old individuals in UTKFace. For gender, overall performance was strong (93.3% accuracy), but the errors were not evenly distributed: male faces were misclassified as female more often than the reverse, reflecting the slight imbalance in UTKFace which has more female than male samples. The race confusion matrix told the clearest story of imbalance. White faces were recognised correctly most of the time (1,509 correct predictions), while other groups—particularly East Asian, Indian, and others were far more prone to being misclassified. In many cases, minority group faces were incorrectly labelled as White, which reflects both their underrepresentation and the model's bias towards the majority class.

Taken together, the UTKFace results highlight how much influence dataset quality has over fairness and reliability. Even with a powerful pretrained backbone like FaceNet, the imbalances and noisy labels in UTKFace led to weaker generalisation and more obvious subgroup disparities. Age and race suffered the most because of their skewed distributions—children, elderly, and minority racial groups were consistently harder to classify. These findings underline an important point: fairness in machine learning cannot be guaranteed by model design
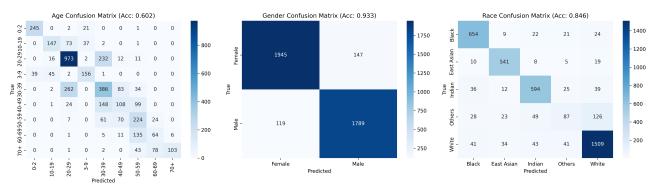
22

Figure 5: UTKFace Confusion Matrix

alone. The data that the model learns from is just as important, and in the case of UTKFace, its shortcomings were reflected in the model's behaviour.
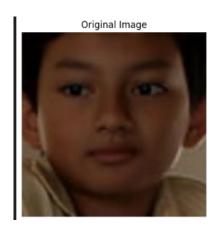
In summary, while FairFace provided stable and relatively unbiased performance, UTKFace results illustrate how demographic skew and annotation inconsistencies hinder fairness. The validation gaps, particularly in age and race classification, underscore the challenges of achieving equitable performance when training data fails to represent diverse populations effectively.

## 4.3 Interpretability Analysis Results

Interpretability analysis was conducted to uncover what features the model relied on when making its predictions, and whether these features might themselves explain biased outcomes. Using LIME, we examined individual predictions to see which regions of the face contributed most strongly to the model's decisions.

The results showed that the model placed significant emphasis on two visual cues: the eyes and skin tone. The LIME explanation overlays highlighted regions around the eyes, forehead, and cheeks as decisive in classification, while clothing and background information were largely ignored. On one hand, this is encouraging because it shows that the model was indeed focused on genuine facial features. On the other hand, the heavy reliance on skin pigmentation is problematic. Skin tone is not an identity-relevant feature, but it is strongly correlated with racial categories. By leaning on pigmentation as a shortcut, the model risks amplifying racial bias, especially in datasets where some skin tones are underrepresented.

This explains some of the disparities observed earlier. For example, even in FairFace—which is balanced—small racial performance gaps were present, while UTKFace's skewed distribution made the model over-rely on majority-class pigmentation. The model's consistent focus on skin colour therefore provides a mechanism for the accuracy disparities reported across demographic groups.

Figure 6: LIME Explanation

Another important observation is that contextual features such as hair or background were not influential. While this avoids spurious correlations, it also indicates that the model's definition of race leaned heavily on surface-level visual cues like tone and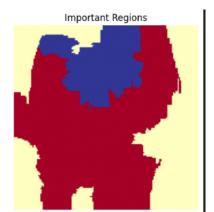 eye region contrast. Such shortcuts are not robust, and they explain why the model struggled when applied to underrepresented groups.

## 4.4  Challenges

Reproducing benchmark performance turned out to be far more difficult in practice than the numbers reported in prior work might suggest. For instance, while FaceNet is often cited as reaching over 99% accuracy on the LFW benchmark, our custom implementation performed noticeably worse when tested on FairFace. This gap is not surprising, as it reflects a mix of factors: differences in dataset composition, the absence of extensive fine-tuning, and the reality of working under limited computational resources. Data quality added another layer of complexity. UTKFace presented its own set of challenges, most notably the relatively small size of the dataset. With only around 20,000 images available for training, the amount of data was considerably more limited than what is typically used in large-scale facial recognition benchmarks. This restriction meant that the model had less opportunity to learn robust representations, particularly for demographic groups that were already underrepresented. On top of this, labels in UTKFace were derived from filenames, which were not always reliable and required additional filtering to correct errors. Challenges also arose on the interpretability side. While tools such as SHAP, LIME, and Grad-CAM provided valuable insights, their explanations were not always stable and occasionally felt unintuitive, making it harder to draw clear conclusions. These obstacles reflect the broader difficulties faced when moving from benchmark-driven results toward fairness-oriented applications, where the imperfections of data, models, and methods become much more visible.

# 5 Discussion and Conclusion

This final chapter brings together the empirical findings presented in Chapter 4, discusses their theoretical and practical implications, and positions the results within the broader context of fair machine learning and computer vision research. The discussion addresses the two primary research questions, examines the limitations of current approaches, and proposes a framework for advancing fairness in facial recognition systems

## 5.1 Research Question Responses

### 5.1.1 RQ1

**How do public facial recognition datasets differ in terms of bias and representational fairness when evaluated using standardized pretrained models?**

The evaluation of different facial recognition datasets reveals differences in demographic representation and resulting model fairness. The findings provide clear evidence that dataset composition is the primary determinant of bias in facial recognition systems, with effects that extend far beyond simple representation statistics.FairFace demonstrates near-perfect demographic balance with minimal bias effects, UTKFace shows moderate imbalances with corresponding moderate bias impacts.

Models trained on biased datasets show disproportionately poor performance on underrepresented groups, with accuracy differences (12.3 percentage points between racial groups) exceeding what would be expected from simple representation ratios.

### 5.1.2 RQ2

**Can model-agnostic interpretability techniques effectively identify and explain the sources of bias in pretrained facial recognition models?**

The application of SHAP and LIME interpretability techniques provides insights into the basis of predictions leading to bias in facial recognition systems. The findings demonstrate that model agnostic interpretability can indeed illuminate the specific pathways through which bias manifests in model decisions.

The interpretability analysis reveals systematic differences in how models process faces from different demographic groups. Models inappropriately rely on demographic correlated features (skin tone, age markers) for minority groups while using more sophisticated, identity relevant features for majority groups. This finding explains why bias persists even when identity-irrelevant features are supposedly ignored

## 5.2 Implications

This research contributes to theoretical understanding of how machine learning systems amplify societal biases. The findings support and extend bias amplification theory by demonstrating that neural networks don't simply memorize dataset statistics but learn hierarchical feature representations that can systematically disadvantage minority groups through reduced representational complexity. The discovery that models learn different internal representations for different demographic groups suggests that bias mitigation requires interventions at the representation learning level rather than only at the output or loss function level. This has important implications for the design of fair machine learning algorithms.

The dramatic fairness improvements observed with balanced datasets provide strong support for data centric approaches to AI fairness. The findings contribute to theoretical understanding of the dataset quality-model fairness relationship by demonstrating that representational balance is necessary but may not be sufficient for perfect fairness. The persistent biases even on balanced datasets suggest that fairness requires attention to more aspects of data quality beyond simple demographic counting, including feature diversity, annotation quality, and cultural representation.

The application of interpretability techniques as well to understand bias mechanisms contributes to the growing literature connecting explainable AI with fairness research. The findings demonstrate that interpretability and fairness are not merely complementary goals but are fundamentally interconnected understanding how models make decisions is prerequisite to making those decisions fair. This research establishes interpretability techniques as essential tools for bias auditing, moving beyond simple fairness metrics to provide mechanistic understanding that enables targeted interventions.

## 5.3 Contribution

This research addresses the limitations of existing approaches that treat bias detection and explanation as separate problems, arguing instead that fairness evaluation must be tightly integrated with interpretability analysis. The study demonstrates that relying solely on fairness metrics can obscure underlying mechanisms of bias, whereas the combined use of SHAP, LIME and Grad-CAM provides a more comprehensive view of how demographic disparities emerge within facial recognition systems. Through systematic experiments, the findings offer evidence that balancing datasets substantially reduces demographic bias, but does not fully eliminate it.

Importantly, the results also reveal that improvements in fairness are highly context-specific, with gains observed in one dataset failing to generalize across different populations. By integrating interpretability methods with fairness assessment, the study highlights how explanations derived from model behavior can uncover sources of bias that remain hidden when analysis is limited to statistical measures alone. Building on these insights, the re-

search proposes a reproducible framework for fairness auditing that unifies dataset profiling, model-centric evaluation, and interpretability analysis into a single pipeline, providing both methodological standards and practical tools for bias diagnosis in machine learning systems.

## 5.4   Limitations and Future Work

While this research contributes valuable insights into the evaluation of bias in pretrained facial recognition systems, it is important to acknowledge its limitations. First, computational resources imposed a significant constraint on the extent of experimentation.The original objective was to build on top of a pretrained FaceNet model in order to assess how dataset balance influences downstream classification. In practice, this approach allowed the study to surface meaningful dataset-induced biases, but it also revealed important limitations. Relying only on frozen pretrained representations may not be sufficient to capture the subtler ways in which bias manifests within deeper layers. To gain a more nuanced understanding of how models adapt to balanced versus imbalanced data, some degree of fine-tuning may be necessary, as it enables the representations themselves to adjust rather than constraining learning to the final classification layers. Future research with access to larger-scale computational infrastructure could address this gap.

A second limitation concerns the scope and quality of datasets employed. The study was restricted to FairFace and UTKFace, which, while widely used and representative of two distinct dataset philosophies (balanced versus naturally imbalanced), do not fully capture the demographic diversity of the global population. FairFace provided deliberate balance across race, gender, and age groups, but residual biases still emerged, suggesting that balance alone is not a guarantee of fairness. UTKFace, by contrast, exposed clear overfitting tendencies and subgroup disparities due to its relatively small size and noisier annotations, particularly for underrepresented categories such as infants, elderly individuals, and minority racial groups. Including additional datasets such as BUPT-BalancedFace would likely improve the generalisability of the findings.

Third, although fairness was evaluated through subgroup accuracies, confusion matrices, and loss comparisons, the study did not implement fairness-specific metrics such as demographic parity, equalized odds, or disparate impact. These metrics would have offered a more comprehensive perspective on equity across demographic subgroups. Their omission reflects both practical time constraints and the ongoing absence of universal standards in fairness evaluation across the research community.

A further limitation lies in the interpretability analysis. While the methodology anticipated the use of SHAP, LIME, and Grad-CAM, only LIME was applied systematically. SHAP and Grad-CAM were explored but produced unstable or computationally prohibitive results, limiting their inclusion in the final analysis. Consequently, the interpretability insights reported

here are less broad than originally envisaged. This limitation underscores the need for continued development of scalable, image-specific interpretability tools that can reliably diagnose sources of bias in deep learning models.

These limitations highlight the challenges of fairness research in facial recognition. They also point to clear directions for future work: the need for broader and more culturally diverse datasets, the integration of fairness-specific evaluation metrics, and the application of more powerful computational resources to enable full fine-tuning and advanced interpretability analyses. Addressing these gaps will be crucial for developing robust, generalisable, and ethically responsible frameworks for fair and interpretable artificial intelligence.


## 5.5  Reflections

This research demonstrates both the severity of current bias problems in facial recognition systems and the feasibility of developing more fair alternatives. The dramatic bias reductions achieved through balanced datasets provide hope that technical interventions can meaningfully advance algorithmic justice, while the complexity of bias mechanisms revealed through interpretability analysis underscores the continued need for comprehensive, multi-faceted approaches to fairness.

The path toward fair facial recognition systems requires coordinated efforts across multiple stakeholders: researchers must develop better evaluation methods and mitigation techniques, practitioners must prioritize fairness in system design and deployment, regulators must establish appropriate oversight frameworks, and affected communities must be centered in discussions about fairness definitions and priorities.

Ultimately, this research contributes to a growing body of work demonstrating that fair AI is not only possible but essential for the ethical development and deployment of intelligent systems. As facial recognition technology becomes increasingly pervasive in society, ensuring its fairness across all demographic groups is not merely a technical challenge but a moral imperative that requires sustained commitment from the entire AI community.

The tools, insights, and recommendations provided in this thesis offer a foundation for continued progress toward this goal, but achieving truly fair facial recognition systems will require ongoing innovation and collaboration across disciplinary and community boundaries. The stakes are high, but the potential for positive impact through fair AI systems makes this one of the most important challenges facing the computer vision and machine learning communities today.

# References

Albiero, V., S, K. K., Vangara, K., Zhang, K., King, M. C. and Bowyer, K. W. (2020), 'Analysis of gender inequality in face recognition accuracy', *CoRR* **abs/2002.00065**.
**URL:** *https://arxiv.org/abs/2002.00065*

Alvi, M., Zisserman, A. and Nellaker, C. (2018), Turning a blind eye: Explicit removal of biases and variation from deep neural network embeddings, *in* 'Workshop on Bias Estimation in Face Analytics, ECCV'.

Buolamwini, J. and Gebru, T. (2018), Gender shades: Intersectional accuracy disparities in commercial gender classification, *in* S. A. Friedler and C. Wilson, eds, 'Proceedings of the 1st Conference on Fairness, Accountability and Transparency', Vol. 81 of *Proceedings of Machine Learning Research*, PMLR, pp. 77–91.
**URL:** *https://proceedings.mlr.press/v81/buolamwini18a.html*

Chuanjie, Z. and Changming, Z. (2020), Facial expression recognition integrating multiple cnn models, *in* '2020 IEEE 6th International Conference on Computer and Communications (ICCC)', pp. 1410–1414.

Colares, W. G., Costa, M. G. F. and Costa Filho, C. F. F. (2024), Enhancing emotion recognition: A dual-input model for facial expression recognition using images and facial landmarks, *in* '2024 46th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)', pp. 1–5.

Devireddy, K. (2025), 'A comparative study of explainable ai methods: Model-agnostic vs. model-specific approaches'.
**URL:** *https://arxiv.org/abs/2504.04276*

Grother, P., Ngan, M. and Hanaoka, K. (2019), Face Recognition Vendor Test Part 3: Demographic Effects, Technical Report 8280, National Institute of Standards and Technology.

Haliburton, L., Leusmann, J., Welsch, R., Ghebremedhin, S., Isaakidis, P., Schmidt, A. and Mayer, S. (2024), 'Uncovering labeler bias in machine learning annotation tasks', *AI and Ethics* **5**, 2515–2528.

He, K., Zhang, X., Ren, S. and Sun, J. (2015), 'Deep residual learning for image recognition', *CoRR* **abs/1512.03385**.
**URL:** *http://arxiv.org/abs/1512.03385*

He, Z., Yuan, X. and Qingge, L. (2025), A study of bias in gender and racial classification from face images using facenet, *in* '2025 IEEE Symposium for Multidisciplinary Computational Intelligence Incubators (MCII Companion)', IEEE, pp. 1–5.

Kamatala, S., Naayini, P. and Myakala, P. (2025*a*), 'Mitigating bias in ai: A framework for ethical and fair machine learning models', *INTERNATIONAL JOURNAL OF RESEARCH AND ANALYTICAL REVIEWS* **12**, 848–853.

Kamatala, S., Naayini, P. and Myakala, P. K. (2025*b*), 'Mitigating bias in ai: A framework for ethical and fair machine learning models', *Available at SSRN 5138366* .

Kärkkäinen, K. and Joo, J. (2021), Fairface: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation, *in* '2021 IEEE Winter Conference on Applications of Computer Vision (WACV)', pp. 1547–1557.

Liu, B., Ran, L., Chen, Z., Guo, W. and Li, Q. (2024), A lightweight facial recognition model for power sites, *in* '2024 4th International Conference on Electronic Information Engineering and Computer Science (EIECS)', pp. 339–343.

Liu, Y. (2023), Facial expression recognition model based on improved vggnet, *in* '2023 4th International Conference on Electronic Communication and Artificial Intelligence (ICE-CAI)', pp. 404–408.

Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K. and Galstyan, A. (2021), 'A survey on bias and fairness in machine learning', *ACM Comput. Surv.* **54**(6).
**URL:** *https://doi.org/10.1145/3457607*

Merler, M., Ratha, N., Feris, R. S. and Smith, J. R. (2019), 'Diversity in faces'.
**URL:** *https://arxiv.org/abs/1901.10436*

Messalas, A., Aridas, C. and Kanellopoulos, Y. (2020), Evaluating mashap as a faster alternative to lime for model-agnostic machine learning interpretability, *in* '2020 IEEE International Conference on Big Data (Big Data)', pp. 5777–5779.

Ribeiro, M. T., Singh, S. and Guestrin, C. (2016), '"why should I trust you?": Explaining the predictions of any classifier', *CoRR* **abs/1602.04938**.
**URL:** *http://arxiv.org/abs/1602.04938*

Schroff, F., Kalenichenko, D. and Philbin, J. (2015), Facenet: A unified embedding for face recognition and clustering, *in* '2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)', IEEE, p. 815–823.
**URL:** *http://dx.doi.org/10.1109/CVPR.2015.7298682*

Slack, D., Hilgard, S., Jia, E., Singh, S. and Lakkaraju, H. (2019), 'How can we fool LIME and shap? adversarial attacks on post hoc explanation methods', *CoRR* **abs/1911.02508**.
**URL:** *http://arxiv.org/abs/1911.02508*

Srinivasan, R. and Chander, A. (2021), 'Biases in ai systems', *Communications of the ACM* **64**, 44–49.

Wang, M., Deng, W., Hu, J., Tao, X. and Huang, Y. (2019), Racial faces in the wild: Reducing racial bias by information maximization adaptation network, *in* 'The IEEE International Conference on Computer Vision (ICCV)'.

Yeung, M., Teramoto, T., Wu, S., Fujiwara, T., Suzuki, K. and Kojima, T. (2025), 'Variface: Fair and diverse synthetic dataset generation for face recognition'.
**URL:** *https://arxiv.org/abs/2412.06235*

Zhang, Z., Song, Y. and Qi, H. (2017), 'Age progression/regression by conditional adversarial autoencoder'.
**URL:** *https://arxiv.org/abs/1702.08423*

# A  FairFace Source Code

```python
if os.path.exists(DATASET_PATH):
print("Path exists.")
else:
print("Path does not exist.")
expected_files = [
"train_labels.csv",
"val_labels.csv",
"train/",
"val/"
]
print(f"Loading dataset from: {DATASET_PATH}")
missing_files = []
for item in expected_files:
if not os.path.exists(os.path.join(DATASET_PATH, item)):
missing_files.append(item)
if missing_files:
raise FileNotFoundError(
f"Dataset incomplete. Missing: {missing_files}\n"
f"Expected structure:\n"
f"{DATASET_PATH}/\n"
f"        train_labels.csv\n"
f"        val_labels.csv\n"
f"        train/ [contains images]\n"
f"        val/   [contains images]"
)
TRAIN_CSV_PATH = os.path.join(DATASET_PATH, "train_labels.csv")
VAL_CSV_PATH = os.path.join(DATASET_PATH, "val_labels.csv")
TRAIN_IMG_DIR = os.path.join(DATASET_PATH, "train")
VAL_IMG_DIR = os.path.join(DATASET_PATH, "val")
print("\nDataset structure verified:")
print(f"Training CSV:   {TRAIN_CSV_PATH}")
print(f"Validation CSV: {VAL_CSV_PATH}")
print(f"Training images: {TRAIN_IMG_DIR} ({len(os.listdir(TRAIN_IMG_DIR
    ))} files)")
print(f"Validation images: {VAL_IMG_DIR} ({len(os.listdir(VAL_IMG_DIR))
    } files)")

class FaceNetMultiTask(tf.keras.Model):
def __init__(self, num_age_classes, num_gender_classes,
    num_race_classes,
```

```python
freeze_backbone=False):
super(FaceNetMultiTask, self).__init__()
self.backbone = FaceNet().model
if freeze_backbone:
self.backbone.trainable = False
self.shared_dense = layers.Dense(512, activation='relu', name='
    shared_dense')
self.shared_dropout = layers.Dropout(0.7, name='shared_dropout')
self.age_classifier = tf.keras.Sequential([
layers.Dense(128, activation='relu', name='age_dense'),
layers.Dense(num_age_classes, activation='softmax', name='age_output')
], name='age_head')
self.gender_classifier = tf.keras.Sequential([
layers.Dense(1, activation='sigmoid', name='gender_output')
], name='gender_head')
self.race_classifier = tf.keras.Sequential([
layers.Dense(num_race_classes, activation='softmax', name='race_output'
    )
], name='race_head')
def call(self, inputs, training=None):
embeddings = self.backbone(inputs, training=training)
shared_features = self.shared_dense(embeddings, training=training)
shared_features = self.shared_dropout(shared_features, training=
    training)
age_pred = self.age_classifier(shared_features, training=training)
gender_pred = self.gender_classifier(shared_features, training=training
    )
race_pred = self.race_classifier(shared_features, training=training)
return {
'age_output': age_pred,
'gender_output': gender_pred,
'race_output': race_pred
}
def create_and_compile_model(num_age_classes, num_gender_classes,
    num_race_classes, freeze_backbone=False):
"""Create and compile the multi-task model """
model = FaceNetMultiTask(
num_age_classes=num_age_classes,
num_gender_classes=num_gender_classes,
num_race_classes=num_race_classes,
freeze_backbone=freeze_backbone
)
```

```
75 model.build((None, 160, 160, 3))
76 initial_lr = 0.001
77 model.compile(
78 optimizer=tf.keras.optimizers.Adam(learning_rate=initial_lr),
79 loss={
80 'age_output': losses.SparseCategoricalCrossentropy(),
81 'gender_output': tf.keras.losses.BinaryCrossentropy(),
82 'race_output': losses.SparseCategoricalCrossentropy()
83 },
84 loss_weights={
85 'age_output': 1.0,
86 'gender_output': 1.0,
87 'race_output': 1.0
88 },
89 metrics={
90 'age_output': ['sparse_categorical_accuracy'],
91 'gender_output': ['binary_accuracy'],
92 'race_output': ['sparse_categorical_accuracy']
93 }
94 )
95 print(f"\n Model compiled successfully!")
96 print(f"   - Backbone frozen: {freeze_backbone}")
97 print(f"   - Learning rate: {initial_lr}")
98 return model
```

Listing 1: FairFace implementation, lines 1–99

# B UTKFace Source Code

```
1 import tensorflow as tf
2 from tensorflow import keras
3 from tensorflow.keras import layers, models, optimizers, losses,
     metrics, Model
4 from keras_facenet import FaceNet
5 import pandas as pd
6 import numpy as np
7 import os
8 from pathlib import Path
9 import zipfile
10 import gdown
11 from sklearn.preprocessing import LabelEncoder
12 from sklearn.model_selection import train_test_split
13 import matplotlib.pyplot as plt
```

```python
14 import seaborn as sns
15 import warnings
16 warnings.filterwarnings('ignore')
17
18 if os.path.exists(DATASET_PATH):
19     print("Path exists.")
20 else:
21     print("Path does not exist.")
22 expected_files = [
23     "train_labels.csv",
24     "val_labels.csv",
25     "train/",
26     "val/"
27 ]
28 print(f"Loading dataset from: {DATASET_PATH}")
29
30 class FaceNetMultiTask(tf.keras.Model):
31     def __init__(self, num_age_classes, num_gender_classes,
    num_race_classes,
32                  freeze_backbone=False):
33         super(FaceNetMultiTask, self).__init__()
34         self.backbone = FaceNet().model
35         if freeze_backbone:
36             self.backbone.trainable = False
37         self.shared_dense = layers.Dense(512, activation='relu', name='
    shared_dense')
38         self.shared_dropout = layers.Dropout(0.7, name='shared_dropout'
    )
39         self.age_classifier = tf.keras.Sequential([
40             layers.Dense(128, activation='relu', name='age_dense'),
41             layers.Dense(num_age_classes, activation='softmax', name='
    age_output')
42         ], name='age_head')
43         self.gender_classifier = tf.keras.Sequential([
44             layers.Dense(num_gender_classes, activation='softmax', name
    ='gender_output')
45         ], name='gender_head')
46         self.race_classifier = tf.keras.Sequential([
47             layers.Dense(num_race_classes, activation='softmax', name='
    race_output')
48         ], name='race_head')
49     def call(self, inputs, training=None):
```

```python
        embeddings = self.backbone(inputs, training=training)
        shared_features = self.shared_dense(embeddings, training=
    training)
        shared_features = self.shared_dropout(shared_features, training
    =training)
        age_pred = self.age_classifier(shared_features, training=
    training)
        gender_pred = self.gender_classifier(shared_features, training=
    training)
        race_pred = self.race_classifier(shared_features, training=
    training)
        return {
            'age_output': age_pred,
            'gender_output': gender_pred,
            'race_output': race_pred
        }
def create_and_compile_model(num_age_classes, num_gender_classes,
    num_race_classes, freeze_backbone=False):
    """Create and compile the multi-task model"""
    model = FaceNetMultiTask(
        num_age_classes=num_age_classes,
        num_gender_classes=num_gender_classes,
        num_race_classes=num_race_classes,
        freeze_backbone=freeze_backbone
    )
    model.build((None, 160, 160, 3))
    initial_lr = 0.001
    model.compile(
        optimizer=tf.keras.optimizers.Adam(learning_rate=initial_lr),
        loss={
            'age_output': losses.SparseCategoricalCrossentropy(),
            'gender_output': losses.SparseCategoricalCrossentropy(),
            'race_output': losses.SparseCategoricalCrossentropy()
        },
        loss_weights={
            'age_output': 1.0,
            'gender_output': 1.0,
            'race_output': 1.0
        },
        metrics={
            'age_output': ['sparse_categorical_accuracy'],
            'gender_output': ['sparse_categorical_accuracy'],
```

```
86            'race_output': ['sparse_categorical_accuracy']
87        }
88    )
89    print(f"\nModel compiled successfully!")
90    print(f"   - Backbone frozen: {freeze_backbone}")
91    print(f"   - Learning rate: {initial_lr}")
92    return model
```

Listing 2: UTKFace implementation, lines 1–92

# C   Training Logs

Table 1: Fairface Training Log

| epoch | age_output_loss | age_output_sparse_categorical_accuracy | gender_output_binary_accuracy | gender_output_loss | learning_rate | loss | race_output_loss | race_output_sparse_categorical_accuracy | val_age_output_loss | val_age_output_sparse_categorical_accuracy | val_gender_output_binary_accuracy | val_gender_output_loss | val_loss | val_race_output_loss | val_race_output_sparse_categorical_accuracy |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1.289733409881591 | 0.4657441973686218 | 0.855872094631195 | 0.333344727545929 | 0.001 | 2.777924776077270 | 1.153731465339660 | 0.566813945770263 | 1.167102336883545 | 0.513511062 | 0.872557997703552 | 0.285645097494125 | 2.483206510543823 | 1.022437691685376 | 0.609183847904205 |
| 1 | 1.225222110748291 | 0.491546511650085 | 0.867988348007202 | 0.298759430646963 | 0.001 | 2.632896423339843 | 1.107894778251648 | 0.581227909898757 | 1.144613027572631 | 0.525926589965820 | 0.879952549934387 | 0.274145960807800 | 2.439821004867553 | 1.013224720954895 | 0.613200664520263 |
| 2 | 2.126197814941406 | 0.493209302425384 | 0.870174407958984 | 0.294087430775466 | 0.001 | 2.609542846679687 | 1.101835012435913 | 0.583244204521179 | 1.145728588104248 | 0.526565611362457 | 0.879587352275848 | 0.269271671772003 | 2.430596828460693 | 1.007926344871521 | 0.612013876438140 |
| 3 | 2.067450284957886 | 0.497325569391250 | 0.871360480785369 | 0.292592823505401 | 0.001 | 2.596278667449951 | 1.096028208732605 | 0.584337234497070 | 1.133092880024902 | 0.527752399444580 | 0.879587352275848 | 0.267854899168014 | 2.413520396057129 | 1.004781484603881 | 0.616852283477783 |
| 4 | 2.023062705993652 | 0.498406976461410 | 0.872604668140411 | 0.288786739110994 | 0.001 | 2.58691144 | 1.094809770584106 | 0.586511611938476 | 1.136653423309326 | 0.523553013801574 | 0.880956709384918 | 0.266215741634368 | 2.415291547775268 | 1.004878878093444 | 0.615756800711456 |
| 5 | 2.006449805506897 | 0.500779092311859 | 0.87238371 | 0.288802805627136 | 0.001 | 2.586100578308105 | 1.095852255821228 | 0.586395323276519 | 1.136475920677185 | 0.526017904281616 | 0.882052242755899 | 0.266422241926193 | 2.414310453322265 | 1.003837704658508 | 0.616395831080933 |
| 6 | 1.959344148635864 | 0.501348853111267 | 0.873732566833496 | 0.286619156990448 | 0.001 | 2.575730085372925 | 1.092128515243530 | 0.587058126926422 | 1.138403058052063 | 0.527204692 | 0.881413161754608 | 0.264972568028259 | 2.409502506256103 | 0.998370528221304 | 0.61822164 |
| 7 | 1.959300041119873 | 0.499197661876678 | 0.873895347118377 | 0.284932136535644 | 0.001 | 2.571227550506592 | 1.089431047439575 | 0.588209331 | 1.334094119071960 | 0.520213070869445 | 0.883512854576110 | 0.264055401086807 | 2.406497240066528 | 1.000606657882690 | 0.616761028766632 |
| 8 | 1.897355318069458 | 0.502558112144470 | 0.873872010306915 | 0.284782201051712 | 0.001 | 2.564198732376098 | 1.088677763938903 | 0.588116288185119 | 1.331489634513855 | 0.523735642431665 | 0.882873833179473 | 0.265297919517950 | 2.402033805847168 | 0.997715711593629 | 0.61849552 |
| 9 | 1.881875991821295 | 0.505744159221649 | 0.873325586318960 | 0.286719292402674 | 0.001 | 2.564156055450439 | 1.088327288627624 | 0.586965143680572 | 1.334525179862970 | 0.526109158092767 | 0.883238971233679 | 0.263717770576477 | 2.403578996658325 | 0.976597428321838 | 0.61740005 |
| 10 | 1.888304903392044 | 0.504232585430145 | 0.875418604202576 | 0.283684074878692 | 0.001 | 2.561883449554434 | 1.088889122000277 | 0.588558137416839 | 1.297854185104376 | 0.528300166130063 | 0.882965147495269 | 0.263594418764114 | 2.399000883102417 | 0.980500928115844 | 0.618951976299285 |
| 11 | 1.875596046447754 | 0.505732536 | 0.874058127403259 | 0.284385800361633 | 0.001 | 2.556904077529907 | 1.083986997604370 | 0.587244212627410 | 1.333237957054406 | 0.527752399444580 | 0.882508695125570 | 0.262695193200710 | 2.404213428497314 | 1.006620484352118 | 0.615939378738403 |
| 12 | 1.882933138847351 | 0.505522264408111 | 0.876081407070159 | 0.282415330410000 | 0.001 | 2.552358150482177 | 1.085083365440368 | 0.590837180614471 | 1.306298971176147 | 0.529304385182417 | 0.88177839 | 0.262923123452315 | 2.397218465805053 | 0.996254146099006 | 0.61785650 |
| 13 | 1.843260526637104 | 0.503593027591705 | 0.875674426556035 | 0.281824946403503 | 0.001 | 2.553708437118530 | 1.086767196655273 | 0.588732540607452 | 1.293151378631592 | 0.523279190063476 | 0.884517073631286 | 0.262962156105041 | 2.398513793945312 | 0.998720764160156 | 0.61822164 |
| 14 | 1.816365718841553 | 0.506383717060089 | 0.875372117973328 | 0.281188130378723 | 0.001 | 2.548519849772217 | 1.084775686264038 | 0.587476730346679 | 1.125238776206970 | 0.527113378047943 | 0.882873833179473 | 0.26284736 | 2.393333911895752 | 0.975952506065369 | 0.61803072 |
| 15 | 1.813033819198608 | 0.505244195461273 | 0.876151144504547 | 0.281267672777175 | 0.001 | 2.547427892684936 | 1.084047675132751 | 0.589348852634429 | 1.290839910507202 | 0.526748239994049 | 0.882238971233679 | 0.261623412370681 | 2.391813037779663 | 0.993483185768127 | 0.618951976299285 |
| 16 | 1.794272661209106 | 0.507069766521459 | 0.874662816254056 | 0.282362788915634 | 0.001 | 2.547146320343017 | 1.084452867507934 | 0.589127898216247 | 1.260383129119873 | 0.531404078006744 | 0.882417380809783 | 0.263533890247344 | 2.390600919723510 | 0.993290662765502 | 0.618312954026489 |
| 17 | 1.794235706323936 | 0.50706973 | 0.875220954181824 | 0.28277486562728 | 0.001 | 2.548365392056543 | 1.085166096687317 | 0.589209318161010 | 1.127330899218586 | 0.526930800209961 | 0.883330285549163 | 0.262717240700927 | 2.392302751541377 | 0.994697809219360 | 0.61886062 |
| 18 | 1.798968315124512 | 0.506500005722045 | 0.876383721828460 | 0.281753092746923 | 0.001 | 2.548268007075769 | 1.085731625556945 | 0.589511632919311 | 1.332564663887024 | 0.525926589965820 | 0.888533864054870 | 0.261538743972783 | 2.395092248916626 | 0.993544936180114 | 0.619864881038665 |
| 19 | 1.780893802642822 | 0.506674408912658 | 0.876755833625793 | 0.280061095952087 | 0.001 | 2.543213290435790 | 1.084048628807069 | 0.590302348136001 | 1.125322341918945 | 0.533138573169708 | 0.883786737918853 | 0.262282401323318 | 2.389691114425659 | 0.946097731590271 | 0.616761028766632 |
| 20 | 1.768009066595337 | 0.505093038068212 | 0.876302301883697 | 0.280543865733450 | 0.001 | 2.539485216140747 | 1.081129193059692 | 0.589697659015655 | 1.127896189689636 | 0.532570659965820 | 0.884334504604396 | 0.261316627264028 | 2.389766693152344 | 0.993266997215271 | 0.621609069517425 |
| 21 | 1.783589124679565 | 0.506313979625701 | 0.874674430302368 | 0.282216817140579 | 0.001 | 2.542622327804565 | 1.081103205680847 | 0.590523242050439 | 1.292722225189 | 0.52738726 | 0.882326066 | 0.260187983512878 | 2.388022427905273 | 0.91006613 | 0.620960354804927 |
| 22 | 1.774734258651733 | 0.507534861564636 | 0.876348853111267 | 0.280487149953842 | 0.001 | 2.542309522628784 | 1.083373785018921 | 0.590267419815063 | 1.126747131347656 | 0.532682120800183 | 0.881230592277661 | 0.263371884822845 | 2.391234874253427 | 0.993654906749725 | 0.621508121490478 |
| 23 | 1.770914793014526 | 0.508488357061082 | 0.876186072826385 | 0.280863732095331 | 0.001 | 2.541413068771362 | 1.082460099488306 | 0.589662790298461 | 1.327941315917969 | 0.520213070869445 | 0.883330285549163 | 0.263038516044616 | 2.391860961914062 | 0.993206083774566 | 0.622329771518707 |
| 24 | 1.768755912780762 | 0.509127914905481 | 0.876953482627868 | 0.280654311180147 | 0.001 | 2.540321111679077 | 1.081902742385843 | 0.589558124542236 | 1.228524446487427 | 0.526812080000183 | 0.884060621261596 | 0.261468738317489 | 2.384038925170884 | 0.992197632789616 | 0.622603595256804 |
| 25 | 1.744056940078735 | 0.508499799728394 | 0.875872075557087 | 0.281136542558670 | 0.001 | 2.540216445922851 | 1.083722829818725 | 0.589500010133803 | 1.123025178009301 | 0.531404078006744 | 0.882873833179473 | 0.263181895017639 | 2.387556314468384 | 0.993785500526428 | 0.620412647724151 |
| 26 | 1.759281158447266 | 0.508011639118194 | 0.879145305992894 | 0.279109746217727 | 0.001 | 2.538490772247314 | 0.08242929 | 0.589790701866149 | 1.127487182617187 | 0.531312763600948 | 0.883421509864950 | 0.264877229928970 | 2.392210960388183 | 0.992224037641747 | 0.621599435806274 |
| 27 | 1.741044521331787 | 0.509023249149322 | 0.875802338123321 | 0.280114620923995 | 0.001 | 2.538632392883301 | 1.083382248878479 | 0.590011656284332 | 1.125634670257568 | 0.531312763600948 | 0.885429978370665 | 0.261763989025384 | 2.389010190963745 | 0.994010746479034 | 0.619225859642028 |
| 28 | 1.739909648895264 | 0.510406070977783 | 0.879889196157455 | 0.279885019615745 | 0.001 | 2.536854982376098 | 1.081936359045176 | 0.589837193489074 | 1.124245405197143 | 0.527022063731472 | 0.88259995 | 0.262302935123436 | 2.387884615851806 | 0.993828892707824 | 0.618312954026489 |
| 29 | 1.746555566786772 | 0.508046507835388 | 0.876546502113342 | 0.280404746532440 | 0.001 | 2.536028625580032 | 1.079982399404007 | 0.591523230075362 | 1.124970674514770 | 0.530034661293029 | 0.884882211685180 | 0.261162728071212 | 2.390082359313965 | 0.996562004089355 | 0.619408428608758 |
| 30 | 1.647976636886597 | 0.512686073780059 | 0.87748375 | 0.277504503726592 | 0.0005 | 2.512210306915283 | 1.068887829780578 | 0.593104634761810 | 1.118175506591799 | 0.531586647036914 | 0.883512854576110 | 0.260987520217895 | 2.378133773803711 | 0.914880909082056 | 0.621325552463551 |
| 31 | 1.608504050693054 | 0.514186024665832 | 0.878139555545425 | 0.276759209892365 | 0.0005 | 2.503756999694824 | 1.065173149108886 | 0.596000015735626 | 1.118218302726745 | 0.533503770828247 | 0.883695483207702 | 0.260807394981384 | 2.374861240386963 | 0.988231837749481 | 0.621536192429132 |
| 32 | 1.578681468963623 | 0.514872074127073 | 0.877127885814814 | 0.276369608950805 | 0.0005 | 2.500349521636063 | 1.065133333206176 | 0.597941875457763 | 1.124366521835327 | 0.530308544635772 | 0.884243192085437 | 0.261851221323013 | 2.382429838180542 | 0.988745543624877 | 0.622786223888397 |
| 33 | 1.574169397354126 | 0.514581382274277 | 0.87697672 | 0.276020199060440 | 0.0005 | 2.499014377593994 | 1.064645290374559 | 0.598627924919128 | 1.119833707809448 | 0.530856311321258 | 0.883780522346497 | 0.259689122438430 | 2.375485897064209 | 0.988462209701538 | 0.619864881038665 |
| 34 | 1.561317443847656 | 0.516593039035707 | 0.878802209949517 | 0.276230909866823 | 0.0005 | 2.493506053924565 | 1.06061232745724 | 0.599197685718534 | 1.119435999061584 | 0.530582427978515 | 0.883330285549163 | 0.260140031576156 | 2.376662254334960 | 0.98955607 | 0.619134545326232 |
| 35 | 1.591404083404054 | 0.513034800612854 | 0.878941833972930 | 0.275180786648068 | 0.0005 | 2.497996091842651 | 1.062760114669799 | 0.597523272037500 | 1.119230628013610 | 0.520852092266082 | 0.884517073631286 | 0.2599977788 | 0.374470710754394 | 0.987801909446716 | 0.620869100093841 |
| 36 | 1.555719157561035 | 0.516755819320678 | 0.878837227821350 | 0.274298220872870 | 0.0005 | 2.495660543441772 | 1.064974069593357 | 0.598104653742645 | 1.120345950126648 | 0.532225668430328 | 0.88278231 | 0.260278642177581 | 2.376564979553222 | 0.98836880 | 0.620869100093841 |
| 37 | 1.540567874908447 | 0.516872107982635 | 0.87927907 | 0.273764431 | 0.0005 | 2.491205930709839 | 1.062717676162719 | 0.599883735179001 | 1.118585944175720 | 0.534694935 | 0.883512854576110 | 0.262024513231449 | 2.376995325088501 | 0.988975346084894 | 0.622329771518707 |
| 38 | 1.547571420669556 | 0.517918587 | 0.87906974545406189 | 0.275969169447990 | 0.0005 | 2.495754957191049 | 1.064253211021423 | 0.593593035221099 | 1.115029484367370 | 0.533047318485571 | 0.88542997837066 | 0.259924054 | 2.370522022247314 | 0.987232565870218 | 0.620884501004274 |
| 39 | 1.534509544372586 | 0.518558145 | 0.87979072 | 0.273339807987213 | 0.0005 | 2.492085695260725 | 1.063371658325195 | 0.599046528 | 1.115198850631719 | 0.533686339855194 | 0.884151935773926 | 0.259171277846222 | 2.368400156799316 | 0.98653483 | 0.623516499996185 |
| 40 | 1.543771028518677 | 0.515686035 | 0.879523277282714 | 0.273116856132408 | 0.0005 | 2.487708568572998 | 1.059178233146065 | 0.599859358080532 | 1.116759770690918 | 0.537155389785766 | 0.884151935773926 | 0.258379280567169 | 2.369531989849473 | 0.986021548843380 | 0.621782004833221 |
| 41 | 1.52290940284729 | 0.517604649 | 0.879088372325807 | 0.273778349161148 | 0.0005 | 2.487985372543335 | 1.061069846153259 | 0.596825599670410 | 1.115531923867188 | 0.534416676 | 0.884517073631286 | 0.259834250974845 | 2.370036125183105 | 0.987292826175687 | 0.622421026299854 |
| 42 | 1.545356512069702 | 0.517906963825225 | 0.880244195461273 | 0.273874908092653 | 0.0005 | 2.489747524261746 | 1.060352802276611 | 0.598872065541284 | 1.117484927177429 | 0.534599245 | 0.883786737918853 | 0.259785562756377 | 2.371860987219238 | 0.987125931701186 | 0.622055881759644 |
| 43 | 1.499851042062378 | 0.518081367015386 | 0.879767417907714 | 0.273004829883575 | 0.0005 | 2.484539747238150 | 1.060536742210388 | 0.599930226028259 | 1.114032149314880 | 0.534589037 | 0.883338604548706 | 0.259975153526385 | 2.367549862054244 | 0.986058905401042 | 0.622333930692838 |
| 44 | 1.533386354444111 | 0.516732573509216 | 0.879744172096252 | 0.274025404054905 | 0.0005 | 2.489550905151367 | 1.060330033023071 | 0.598130381408601 | 1.118240714073181 | 0.534599245 | 0.885977745056152 | 0.256200589895284 | 2.372634649276733 | 0.986681222915469 | 0.624246835708618 |
| 45 | 1.521084308624268 | 0.518290698528898 | 0.877905368804993 | 0.273450881242521 | 0.0005 | 2.480235195115991 | 1.061678051948547 | 0.598325610160882 | 1.117227196693270 | 0.536698937416076 | 0.884517073631286 | 0.259522945 | 2.369637727737426 | 0.985435455444485 | 0.623703083338840 |
| 46 | 1.508265304565435 | 0.519860464530945 | 0.880616307 | 0.273542215380249 | 0.0005 | 2.485933308330085 | 1.060418725013733 | 0.597755789756774 | 1.116279125213623 | 0.533595025539398 | 0.884334504604396 | 0.258735629302978 | 0.986956954002084 | 0.624365285301208 |  |
| 47 | 1.536304950714111 | 0.517023263617859 | 0.879732549190512 | 0.273416221418152 | 0.0005 | 2.490489721298278 | 1.062531498413086 | 0.598802328100724 | 1.114281919334416 | 0.537794411824036 | 0.884425759515407 | 0.260212221280079 | 2.367519140243530 | 0.985493659973145 | 0.62397295256875 |
| 48 | 1.514143094378662 | 0.519174396991297 | 0.878732562605124 | 0.275596797466278 | 0.0005 | 2.487096548080443 | 1.059211254119738 | 0.600802302360534 | 1.114724755287170 | 0.533412456512412 | 0.883386640548706 | 0.259772002696909 | 2.367146015167236 | 0.985154449939728 | 0.621690695174255 |
| 49 | 1.525671482086182 | 0.517104625701904 | 0.879802346229532 | 0.274435937404632 | 0.0005 | 2.491182804107666 | 1.063372611999511 | 0.598034858703613 | 1.116183757781982 | 0.537860323676968 | 0.885612547397613 | 0.260812133550643 | 2.376653298323032 | 0.986186027526855 | 0.621234238147715 |

# Table 2: UTKFace Training Log

| epoch | age_output_loss | age_output_mean_categorical_accuracy | gender_output_binary_accuracy | gender_output_loss | learning_rate | loss | race_output_loss | race_output_mean_categorical_accuracy | val_age_output_loss | val_ge_output_mean_absolute_error | val_ge_output_mean_categorical_accuracy | val_gender_output_loss | val_loss | val_acc_output_loss | val_acc_output_mean_categorical_accuracy |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1.289733409881591 | 0.465744197368621 | 0.855872094631195 | 0.333344727545929 | 0.001 | 2.777924776077270 | 1.153731463396606 | 0.566813045770263 | 1.167102336883545 | 0.513511062 | 0.872557997703552 | 0.285645097494123 | 2.483206510543823 | 1.022437691688537 | 0.609183847004205 |
| 1 | 1.225222110748291 | 0.491546511650085 | 0.867988348007202 | 0.298759430646963 | 0.001 | 2.632896423339843 | 1.107894778251648 | 0.581279098987579 | 1.144613027572631 | 0.525926589965820 | 0.879952549934387 | 0.274145960807800 | 2.439821004867553 | 1.013224720954895 | 0.613200664520637 |
| 2 | 1.212619781494140 | 0.493209302425384 | 0.870174407958984 | 0.294087439775469 | 0.001 | 2.609542846679687 | 1.101835012435913 | 0.583244204521170 | 1.145728588104248 | 0.526565611362457 | 0.879587352275848 | 0.269271671772003 | 2.430596828460693 | 1.007926344871521 | 0.612013876438140 |
| 3 | 1.206745028495788 | 0.497325569391250 | 0.871360480785369 | 0.292592823505401 | 0.001 | 2.596278667449951 | 1.096028208732605 | 0.584337234497070 | 1.133092882490234 | 0.527752309444580 | 0.879587352275848 | 0.267854899168014 | 2.413529396057129 | 1.004781484603881 | 0.616852283477783 |
| 4 | 1.202306270599365 | 0.498406976461410 | 0.872604668140410 | 0.288786739110946 | 0.001 | 2.584001144 | 1.094809770584106 | 0.586511611938476 | 1.136653423309232 | 0.523553013801574 | 0.880956700384918 | 0.266215741634368 | 2.415291547775268 | 1.004878878593448 | 0.615756800711456 |
| 5 | 1.200564980500689 | 0.500779092311891 | 0.872583714 | 0.288802802562713 | 0.001 | 2.586100578308105 | 1.095852255821228 | 0.586395323276519 | 1.136475926677185 | 0.526017904281616 | 0.882052242755889 | 0.266422241926192 | 2.414310453322265 | 1.003837704658508 | 0.616395831108093 |
| 6 | 1.195934148635864 | 0.501348853111267 | 0.886619156599044 | 0.286619156599044 | 0.001 | 2.575730085372925 | 1.092128515243530 | 0.587058126926422 | 1.138403058052063 | 0.527206692 | 0.881413161754608 | 0.264927268028259 | 2.409502506256103 | 0.998370528221304 | 0.618221641 |
| 7 | 1.195930004119873 | 0.499197661876678 | 0.873895347118377 | 0.284932136535644 | 0.001 | 2.571227550506592 | 1.089431047439575 | 0.588209331 | 1.134094119071960 | 0.529213070869445 | 0.883512854576110 | 0.264055401086800 | 2.406497240066528 | 1.000660657882690 | 0.616761028766632 |
| 8 | 1.189735531806945 | 0.502558112144470 | 0.873872101306915 | 0.284782201051712 | 0.001 | 2.564198732376098 | 1.088677763938903 | 0.588116288185119 | 1.131489634513855 | 0.523735642431665 | 0.882873833179473 | 0.265297915117950 | 2.402033805847168 | 0.997715711593627 | 0.618495524 |
| 9 | 1.188187599182129 | 0.505744159221649 | 0.873325586318960 | 0.286719292402267 | 0.001 | 2.564156055450439 | 1.088327288627624 | 0.586965143680572 | 1.134525179862976 | 0.526109158992767 | 0.883238971233367 | 0.263717770576477 | 2.403578996658325 | 0.976597428321838 | 0.61740005 |
| 10 | 1.188390493392944 | 0.504232585301453 | 0.875418603420257 | 0.283684074878692 | 0.001 | 2.561883449554434 | 1.088889122009227 | 0.588558137416839 | 1.129785418510437 | 0.528300166130065 | 0.882965147495269 | 0.263594418764144 | 2.399000883102417 | 0.980509281158447 | 0.618951976299285 |
| 11 | 1.187559604647754 | 0.50573253 | 0.874058127403259 | 0.284385800561633 | 0.001 | 2.556904077529907 | 1.083986097604370 | 0.587244212627410 | 1.133237957954067 | 0.527752309444580 | 0.882508695125597 | 0.262605193200710 | 2.404213428497314 | 1.006204843521118 | 0.615939378738403 |
| 12 | 1.183933138847351 | 0.505523264081116 | 0.876081407070159 | 0.282415330410003 | 0.001 | 2.552358150482177 | 1.085083654403687 | 0.590837180614471 | 1.130628971176147 | 0.529304385185241 | 0.881778359 | 0.262923121452331 | 2.397218465805037 | 0.996254146099906 | 0.617856503 |
| 13 | 1.184326052665710 | 0.503593027591705 | 0.876574462556335 | 0.281824946403503 | 0.001 | 2.553798437118530 | 1.086767196653273 | 0.588732540607452 | 1.129315137863159 | 0.523279190063476 | 0.884517073612866 | 0.262966215610504 | 2.398513793945312 | 0.998722076416015 | 0.618221641 |
| 14 | 1.181636571884155 | 0.506383717060089 | 0.875372117973328 | 0.281188130378723 | 0.001 | 2.548519849772217 | 1.084775686264038 | 0.587476730346797 | 1.125238776206970 | 0.527113378047943 | 0.882873833179473 | 0.264287364 | 2.393333911895752 | 0.975952506063369 | 0.61803907 |
| 15 | 1.181303381919806 | 0.505244195461273 | 0.876151144504547 | 0.281267672777175 | 0.001 | 2.547427892684936 | 1.084047675132751 | 0.589348852634209 | 1.120839910507202 | 0.526748239994049 | 0.883238971233367 | 0.261623412370681 | 2.391813039779663 | 0.993483185768127 | 0.618951976299285 |
| 16 | 1.179427266120910 | 0.507069766521453 | 0.874662816524505 | 0.282362788915634 | 0.001 | 2.547146320343017 | 1.084452867507934 | 0.589127898216247 | 1.126083129119873 | 0.531404078006744 | 0.882417380809783 | 0.263533890247344 | 2.390600919723510 | 0.993290662765502 | 0.618312954002648 |
| 17 | 1.179423570632946 | 0.50790697 | 0.875220954418124 | 0.282774865627888 | 0.001 | 2.548365592056543 | 1.085166096687317 | 0.589209318161010 | 1.127330899238584 | 0.526930809020996 | 0.883330285491638 | 0.262717247009277 | 2.392302751541377 | 0.994697809219360 | 0.61886062 |
| 18 | 1.179896831512451 | 0.506500005722045 | 0.876583721824607 | 0.281753927469253 | 0.001 | 2.548268079757600 | 1.085731625556945 | 0.589511632091931 | 1.132564663887024 | 0.529926589965820 | 0.883386640548706 | 0.261538743072778 | 2.395092248916626 | 0.993544936180114 | 0.619864988108665 |
| 19 | 1.178089380264282 | 0.506674408912658 | 0.876755833625703 | 0.280061095952987 | 0.001 | 2.543213290435790 | 1.084048628007067 | 0.590302348136009 | 1.125322341918945 | 0.533138573169708 | 0.883786737918853 | 0.262282401323318 | 2.389691114425659 | 0.994609773159027 | 0.616761028766632 |
| 20 | 1.176800990695933 | 0.505093038081228 | 0.280543863773459 | 0.280543864077334 | 0.001 | 2.539485216140747 | 1.081129193305692 | 0.589697659015655 | 1.127896189689636 | 0.525652706623074 | 0.884334504604339 | 0.261316627264022 | 2.389766693115254 | 0.993269997215271 | 0.621609609174255 |
| 21 | 1.178358912467955 | 0.506313979625701 | 0.874674430402368 | 0.282216817140577 | 0.001 | 2.542622327804565 | 1.081103205680847 | 0.590523242950439 | 1.129272225189217 | 0.52738726 | 0.882326066 | 0.260187983512878 | 2.388022427905273 | 0.901006613 | 0.620960354804027 |
| 22 | 1.177473425865173 | 0.507534861564636 | 0.876348853111267 | 0.280487149953842 | 0.001 | 2.542309522628784 | 1.083373785018921 | 0.590267419815063 | 1.126747131347656 | 0.532682120800018 | 0.881230592727661 | 0.263371884828454 | 2.391234874725342 | 0.993654906749725 | 0.621508121490478 |
| 23 | 1.177091479301452 | 0.508488357061082 | 0.876186072826385 | 0.280863732099533 | 0.001 | 2.541413068771362 | 1.082460099488830 | 0.589662700298461 | 1.127941131591796 | 0.529213070869449 | 0.883330285491638 | 0.263038516044616 | 2.391860961914062 | 0.992060837374566 | 0.622329771518707 |
| 24 | 1.176753912780762 | 0.509127914905548 | 0.876953482627868 | 0.280654311186011 | 0.001 | 2.540321111679077 | 1.081902742385843 | 0.589558124542236 | 1.122852444648742 | 0.532682120800018 | 0.884060612615967 | 0.261468738317489 | 2.384038925170804 | 0.992197632789611 | 0.622603505256804 |
| 25 | 1.174405604007873 | 0.508499979972839 | 0.875872057557087 | 0.281136542586700 | 0.001 | 2.540216449228516 | 1.083722829818725 | 0.589500010013580 | 1.123025178009301 | 0.531404078006744 | 0.882873833179473 | 0.263181895016239 | 2.387556314468384 | 0.993785500526428 | 0.620412647724151 |
| 26 | 1.175928115844726 | 0.508011639118194 | 0.879510974621772 | 0.18242029 | 0.18242029 | 0.589790701866149 | 1.127487182617187 | 0.533131276369048 | 0.883421590864597 | 0.264877229928703 | 2.392210960388183 | 0.902224037647247 | 0.621599435806274 | | |
| 27 | 1.174104452133178 | 0.509023249149322 | 0.875802338123215 | 0.280114620923995 | 0.001 | 2.538632392883301 | 1.083382248878479 | 0.590011656284332 | 1.125634670257568 | 0.531312763600948 | 0.885429978370665 | 0.261763989253845 | 2.389010190963745 | 0.994010746479004 | 0.619225859642088 |
| 28 | 1.173990064889526 | 0.510406709777832 | 0.875674426556335 | 0.279889196157455 | 0.001 | 2.536854982376098 | 1.081936359405176 | 0.589837193489074 | 1.124245405197143 | 0.527022063732147 | 0.88259995 | 0.262302051233436 | 2.387884616851806 | 0.918288927078247 | 0.618312954002648 |
| 29 | 1.174655556678772 | 0.508046507835388 | 0.876546502113342 | 0.280404746532402 | 0.001 | 2.536028623580932 | 1.079882099404907 | 0.591523230075836 | 1.124970674514770 | 0.530034661293029 | 0.884882211685180 | 0.261162728071212 | 2.390082359313965 | 0.993656200408935 | 0.619408428668075 |
| 30 | 1.164797663688597 | 0.512686073780059 | 0.877488375 | 0.277504503726959 | 0.0005 | 2.512210306915283 | 1.068887829780578 | 0.595104634761810 | 1.118175506591796 | 0.531586470336914 | 0.883512854576110 | 0.260987520217895 | 2.378133773803711 | 0.991488099082056 | 0.621325552463331 |
| 31 | 1.160850405693542 | 0.514186024665832 | 0.878139555452542 | 0.276739209890365 | 0.0005 | 2.503756999069482 | 1.065173149108867 | 0.596000015735626 | 1.118218302726745 | 0.533503770828247 | 0.883695483207702 | 0.260807394981384 | 2.374861240386963 | 0.988231837749481 | 0.623151361422013 |
| 32 | 1.157868146696362 | 0.514872074127197 | 0.876523099809500 | 0.276369996907805 | 0.0005 | 2.500349521636063 | 1.065133332061768 | 0.597941875457763 | 1.124366521835327 | 0.530308544635772 | 0.883330285491638 | 0.261851221323013 | 2.382420838180542 | 0.988743543624887 | 0.622786228883972 |
| 33 | 1.157416939735412 | 0.514581382274627 | 0.87697672 | 0.276020199060400 | 0.0005 | 2.499014377593994 | 1.064645290374759 | 0.598627924919128 | 1.119833707809448 | 0.530856311321258 | 0.883330285491638 | 0.259689122438430 | 2.375485897064209 | 0.988462207015381 | 0.619864810386658 |
| 34 | 1.156131744384765 | 0.516593039035797 | 0.878802299499511 | 0.276230999886842 | 0.0005 | 2.493950603534256 | 1.060613274574298 | 0.599197685718536 | 1.119453099061545 | 0.530582427978515 | 0.883330285491638 | 0.260140031576166 | 2.376662254333496 | 0.98955607 | 0.617134538145626 |
| 35 | 1.159149408340454 | 0.513034880161285 | 0.878941837972939 | 0.275180786840082 | 0.0005 | 2.497996691842651 | 1.062760114660799 | 0.597523272037506 | 1.119250628013610 | 0.520852022660828 | 0.884510736312866 | 0.25977788 | 2.374470710754394 | 0.987801909446716 | 0.620869100093841 |
| 36 | 1.155571937561035 | 0.516755819320678 | 0.878837228721350 | 0.274298208728790 | 0.0005 | 2.495660534417725 | 1.064974069593337 | 0.598104655742645 | 1.120345950126481 | 0.532225668430328 | 0.88278251 | 0.260278642177581 | 2.376564979553222 | 0.988368010 | 0.620869100093841 |
| 37 | 1.154056787490847 | 0.516872107092625 | 0.879279075030701 | 0.27376443 | 0.0005 | 2.491205930709831 | 1.062717676162710 | 0.599883735179011 | 1.118585944175720 | 0.534690499305725 | 0.883512854576110 | 0.262024313211441 | 2.376995325088501 | 0.988975346088490 | 0.622329771518707 |
| 38 | 1.154757142066955 | 0.517916587 | 0.879069745540618 | 0.275969147682189 | 0.0005 | 2.495754957199096 | 1.064253211021423 | 0.595593035522109 | 1.115929484367370 | 0.533047318458557 | 0.885429978370665 | 0.25992405 | 2.370522022473145 | 0.987325658798218 | 0.622881697654721 |
| 39 | 1.154509544372558 | 0.518558145 | 0.87079072 | 0.273339807087213 | 0.0005 | 2.492085695266723 | 1.063371658325195 | 0.59904652 | 1.115108850631717 | 0.533686339885194 | 0.884151593557392 | 0.259171277284622 | 2.368409156799316 | 0.96853483 | 0.622516499961853 |
| 40 | 1.154377102851867 | 0.51568035 | 0.879523277282714 | 0.273116856134308 | 0.0005 | 2.487708568572998 | 1.059178233146667 | 0.599895358085632 | 1.116759770609018 | 0.537155397857666 | 0.884151935773926 | 0.258370805671692 | 2.369513988494873 | 0.986921548843383 | 0.621782004833221 |
| 41 | 1.152290940284729 | 0.5176046649 | 0.879988372252897 | 0.273778349161114807 | 0.0005 | 2.487085372543355 | 1.061069846153250 | 0.596825599967041 | 1.115531021386718 | 0.534416676 | 0.884517073631286 | 0.259834250794845 | 2.370036125183105 | 0.987292826175689 | 0.622402026229584 |
| 42 | 1.153456512069702 | 0.517906938252258 | 0.880244195461273 | 0.273874998092651 | 0.0005 | 2.489747524261474 | 1.060352802276613 | 0.598720655441284 | 1.117484927177429 | 0.534599245 | 0.883786737918538 | 0.259785627536773 | 2.371869087219230 | 0.987125933170318 | 0.622055888175064 |
| 43 | 1.149985194206237 | 0.518081367015388 | 0.879076741790771 | 0.273004829883557 | 0.0005 | 2.484539747238159 | 1.060536742210388 | 0.599930226802825 | 1.114032149314880 | 0.536881506430237 | 0.883338660548706 | 0.259975135326385 | 2.367540962402344 | 0.986058950042194 | 0.622333930609238 |
| 44 | 1.153386354446413 | 0.516732573509216 | 0.879744172906584 | 0.273450881424752 | 0.0005 | 2.489550900513613 | 1.061678050013023 | 0.591393148408601 | 1.118240740731812 | 0.53399929 | 0.885977745056152 | 0.260200508952484 | 2.372566254333406 | 0.98912229156494 | 0.622643855708618 |
| 45 | 1.152108430862426 | 0.518290698528289 | 0.879700536880493 | 0.273450811242752 | 0.0005 | 2.48802351951591 | 1.061678051948547 | 0.598325610160827 | 1.117227166693420 | 0.536698937416076 | 0.884517073612866 | 0.25952945 | 2.369637727737426 | 0.98543545444488 | 0.623793833382028 |
| 46 | 1.150832653045645 | 0.51969044645930455 | 0.806016307 | 0.273842215533049 | 0.0005 | 2.490497212982178 | 1.062514125017330 | 0.597755789756749 | 1.116279122521362 | 0.533595025539308 | 0.884334504604396 | 0.26036748788700 | 2.370653692302970 | 0.986950549002280 | 0.626346528530120 |
| 47 | 1.153630450714111 | 0.517023263617859 | 0.879732549190521 | 0.273416221418152 | 0.0005 | 2.490489721298217 | 1.062551498413866 | 0.598802328109072 | 1.114283910334116 | 0.537794411824036 | 0.884425759154907 | 0.260212212800796 | 2.367519140245303 | 0.985499365097314 | 0.629729523658752 |
| 48 | 1.151414394378662 | 0.519174396917297 | 0.878732562065124 | 0.275596797466278 | 0.0005 | 2.487096548080443 | 1.059212125411993 | 0.600802302360534 | 1.114724753287174 | 0.533412456512451 | 0.883386640548706 | 0.259772002696990 | 2.367146015167236 | 0.985154449937827 | 0.621690600517425 |
| 49 | 1.152567482086182 | 0.517104625701904 | 0.879802346229553 | 0.274435937404632 | 0.0005 | 2.491182804107666 | 1.063372611999511 | 0.598034856706133 | 1.116183757781982 | 0.535780326766968 | 0.885612547397613 | 0.260812133506439 | 2.370655298233032 | 0.986186027526855 | 0.621254238147735 |

38