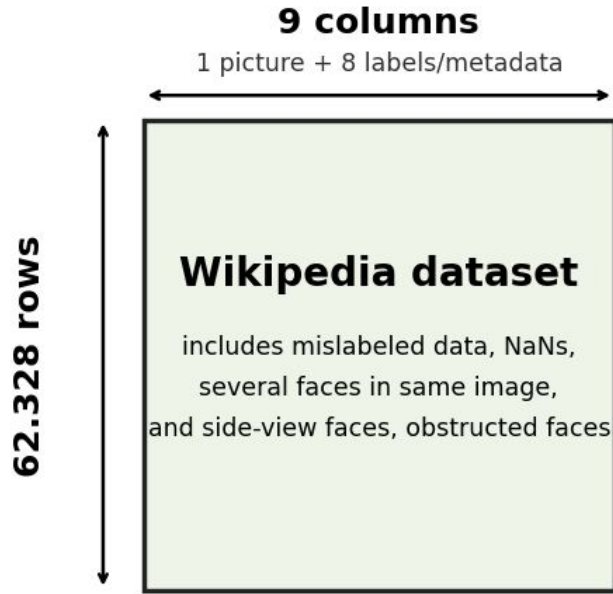


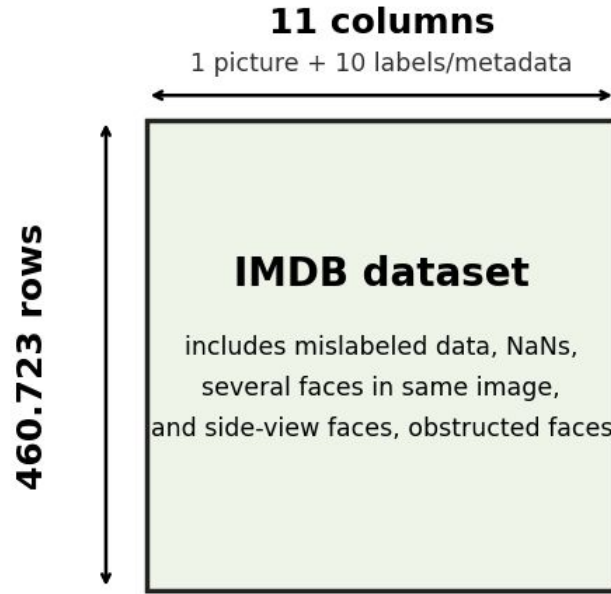
Age & Gender Prediction Using CNN's on Face Images

by Victoria, Jona, David,
Natalie

Data description



labels / metadata:
age, photo_taken, gender, name, face_location, face_score,
second_face_score, birth_date



labels / metadata:
age, photo_taken, gender, name, face_location, face_score,
second_face_score, birth_date

Targets:

Gender

Age

General Cuts

- face_score > 1
 - second_face_score is NaN
 - Age between 0 and 100
 - Gender is **not** NaN
- ~460k to ~170k images

Face-score gradient examples

Face score < 1
score=-inf



Face score > 1
score=2.00



Face score < 1
score=0.73



Face score > 1
score=3.00



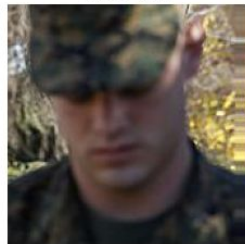
Face score < 1
score=0.79



Face score > 1
score=4.00



Face score < 1
score=0.84



Face score > 1
score=5.00



Face score < 1
score=0.90



Face score > 1
score=6.00



Face score < 1
score=0.95



Face score > 1
score=7.00



Clustering

“Research Questions”

2 Models:

- VGG16
 - Do we find an **underlying structure**?
- Age-Supervised Autoencoder (SAE)
 - Can we **enforce clustering by age**?

General questions:

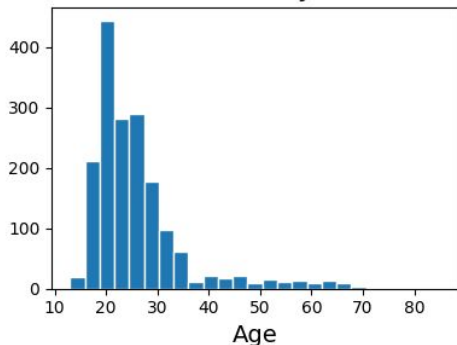
- Influence of the gender?
- Source bias (IMDB vs. Wiki)?

VGG16

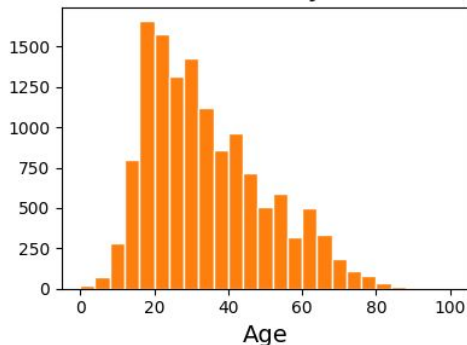
VGG16 - K-Means Clustering

Age-Distribution per Cluster (Sorted by $\bar{\text{Age}}$)

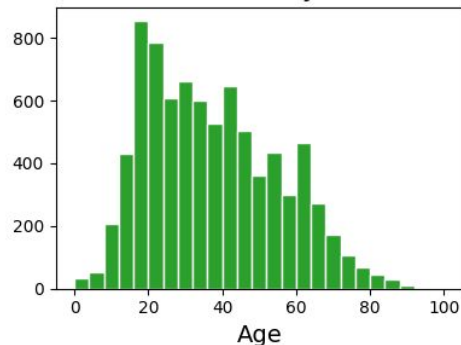
Cluster no. 5, $\bar{\text{Age}}=26$ yrs, $\sigma=10.0$



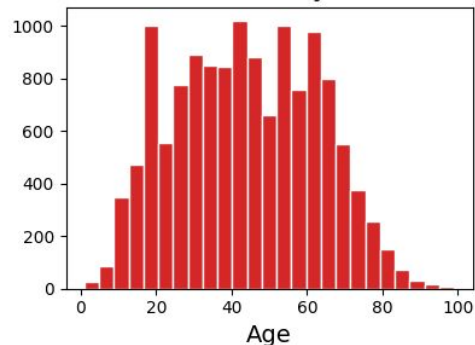
Cluster no. 1, $\bar{\text{Age}}=34$ yrs, $\sigma=16.1$



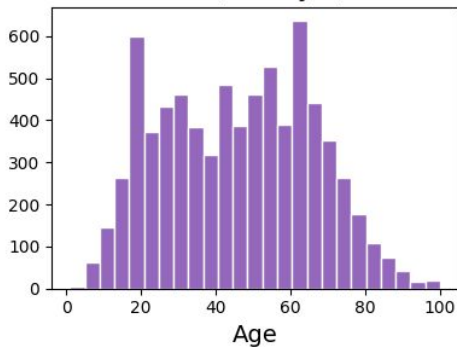
Cluster no. 4, $\bar{\text{Age}}=37$ yrs, $\sigma=17.9$



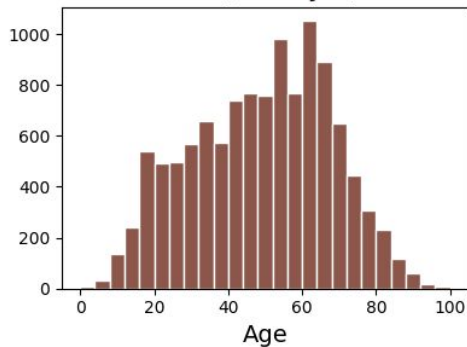
Cluster no. 6, $\bar{\text{Age}}=44$ yrs, $\sigma=18.7$



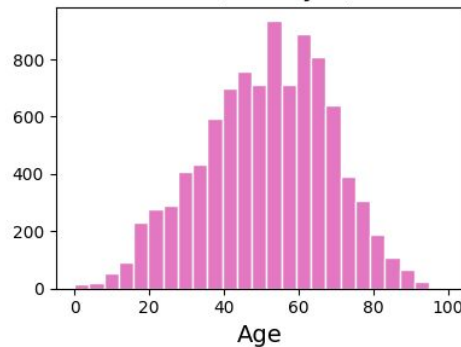
Cluster no. 3, $\bar{\text{Age}}=46$ yrs, $\sigma=20.3$



Cluster no. 0, $\bar{\text{Age}}=48$ yrs, $\sigma=18.8$

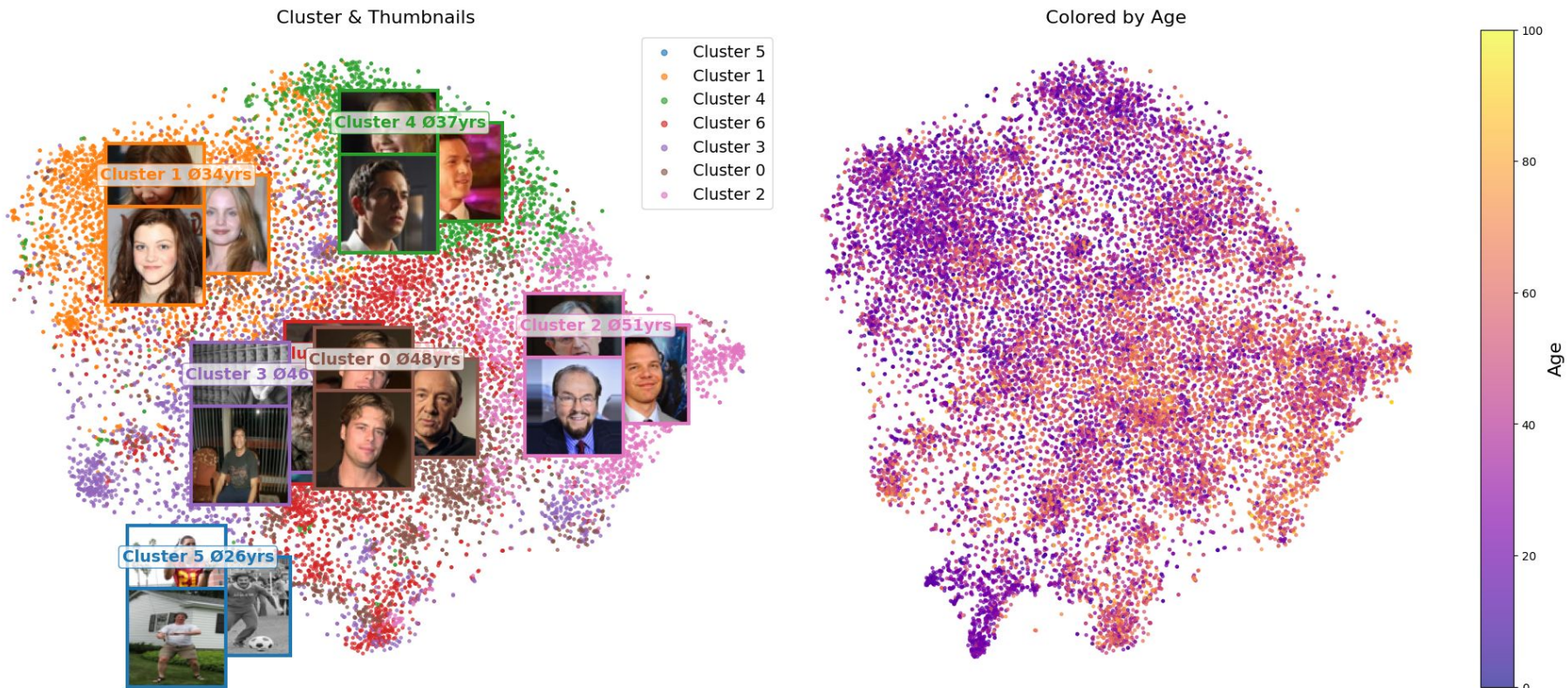


Cluster no. 2, $\bar{\text{Age}}=51$ yrs, $\sigma=17.3$



VGG16 (K-Means) - Visualization by t-SNE

VGG16 fc2 Features - t-SNE Visualization

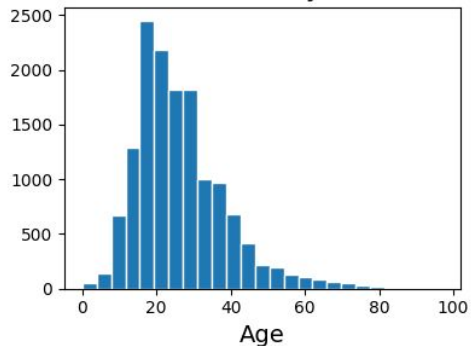


Age-Supervised Autoencoder (SAE)

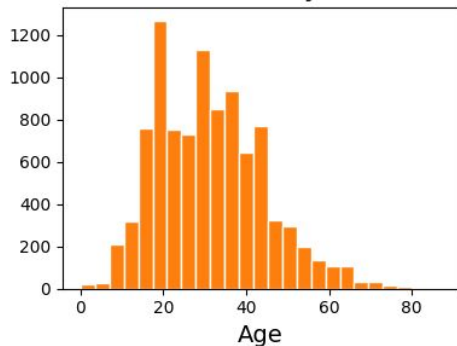
SAE - Age-Distribution of Clusters

Age-Distribution per Cluster (Sorted by \emptyset -Age)

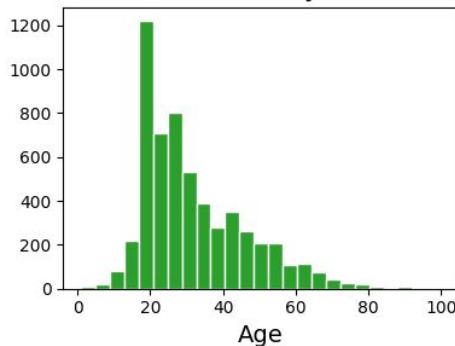
Cluster no. 3, \emptyset 27yrs, σ =12.3



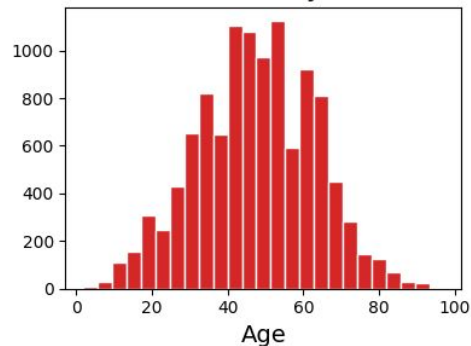
Cluster no. 0, \emptyset 31yrs, σ =13.1



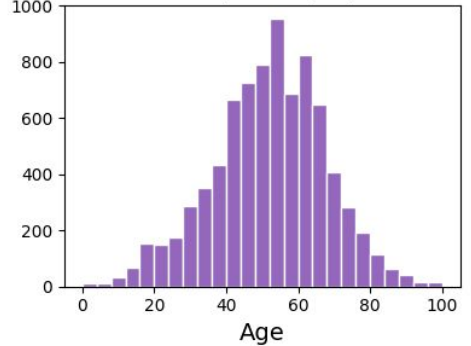
Cluster no. 2, \emptyset 32yrs, σ =14.5



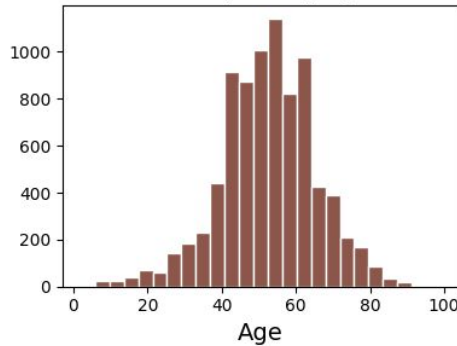
Cluster no. 4, \emptyset 47yrs, σ =15.5



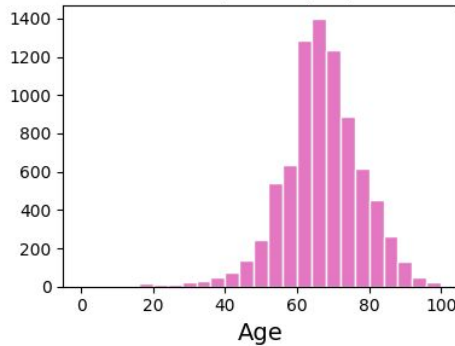
Cluster no. 1, \emptyset 51yrs, σ =15.7



Cluster no. 5, \emptyset 53yrs, σ =13.0



Cluster no. 6, \emptyset 66yrs, σ =11.3

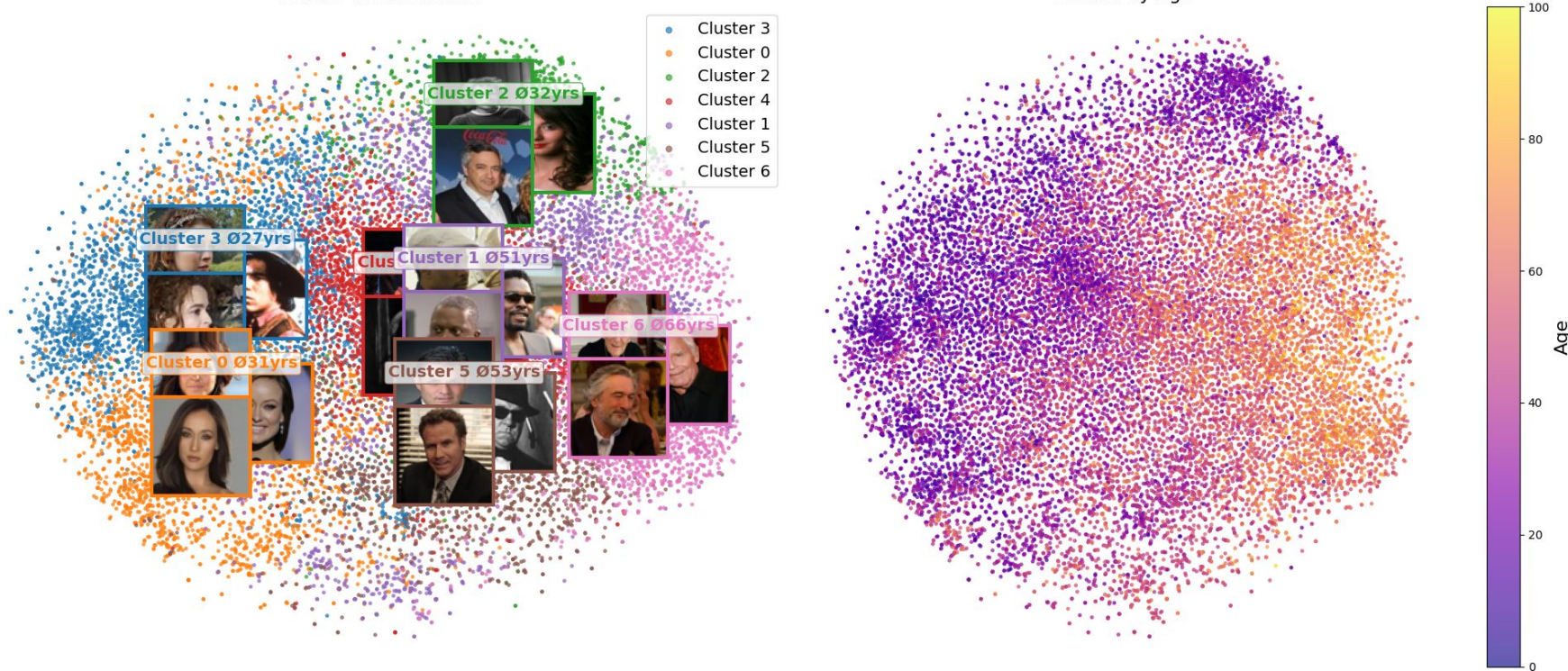


SAE - Visualization by t-SNE

t-SNE - Age-Supervised AE (weight=0.2)

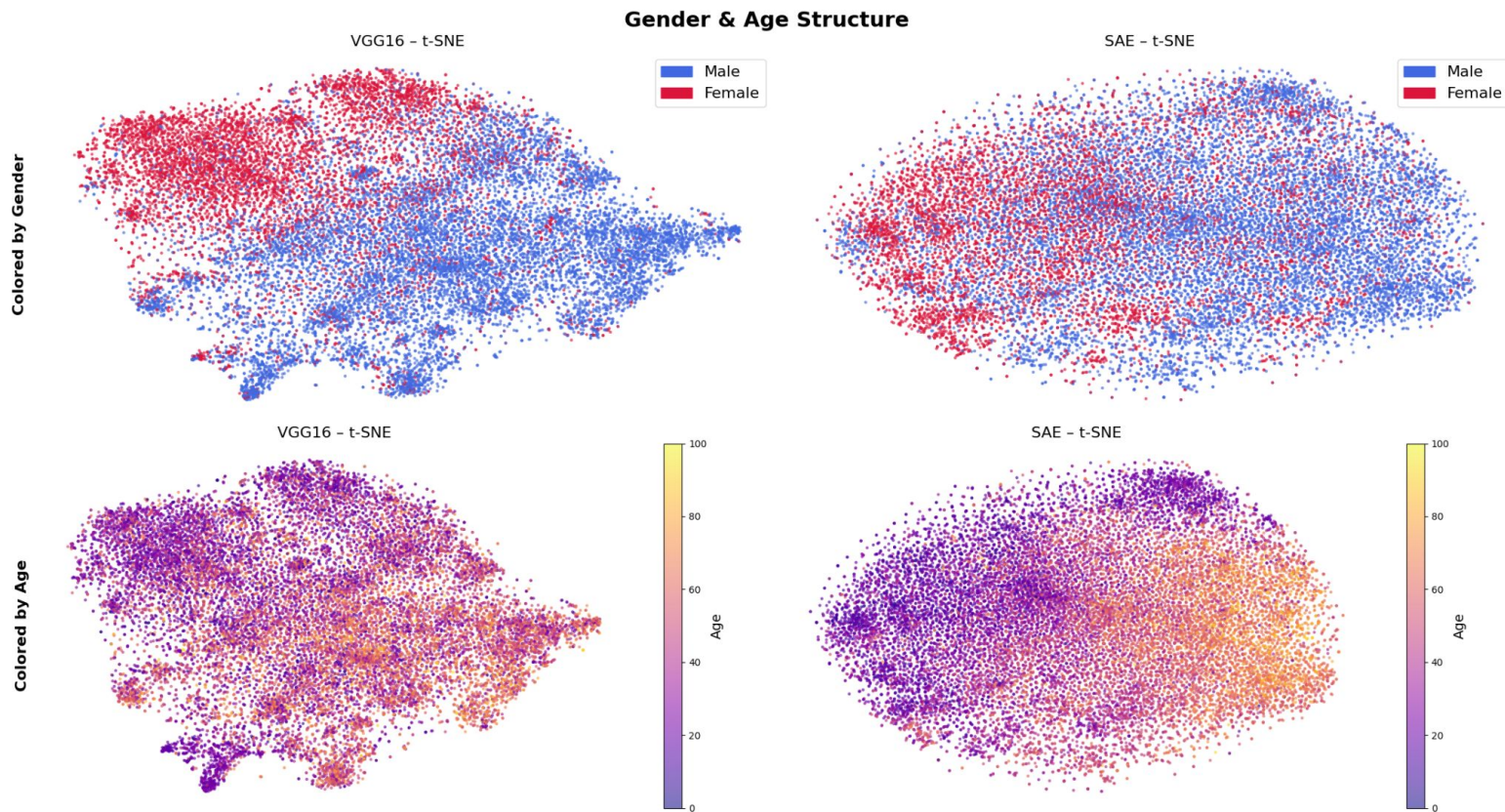
Cluster & Thumbnails

Colored by Age



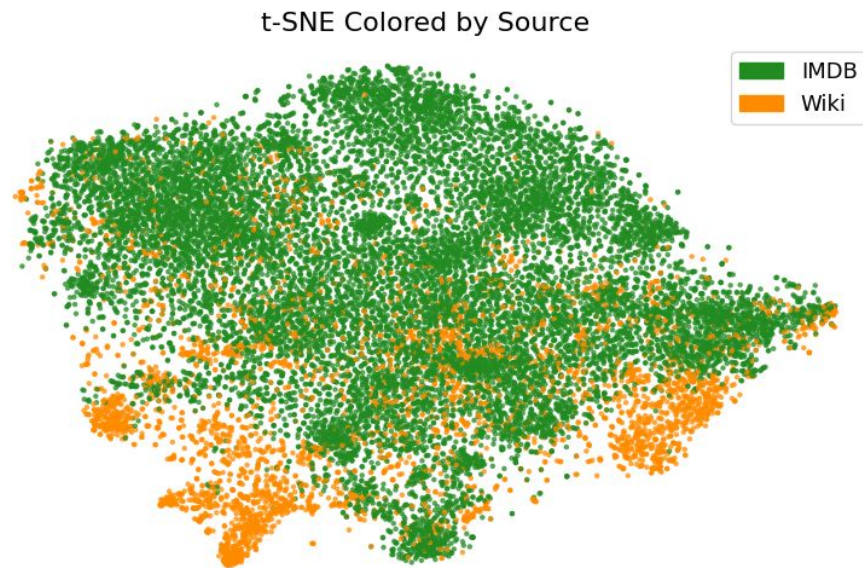
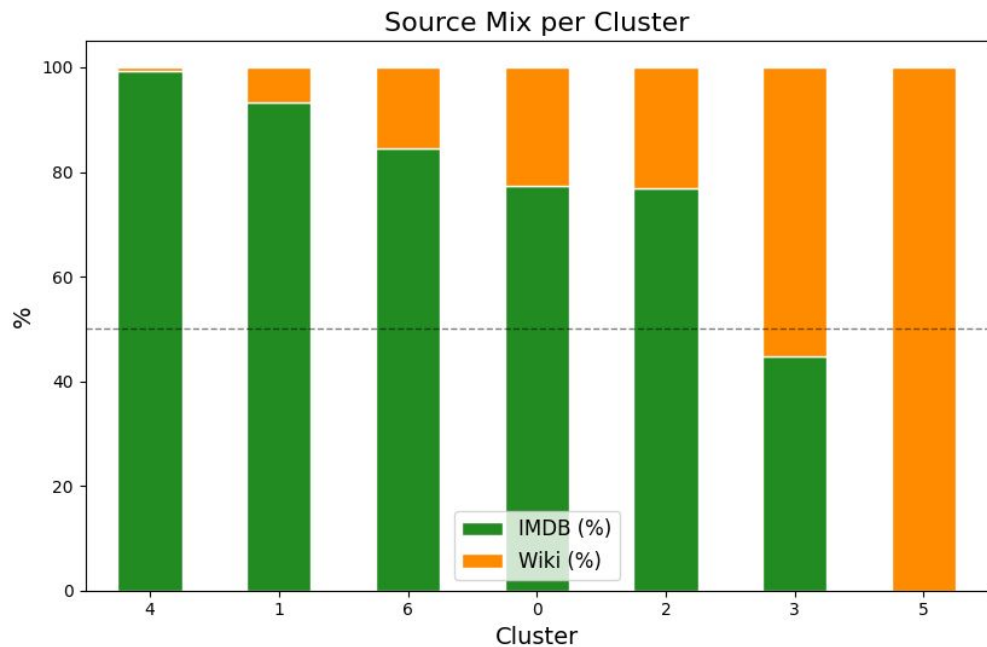
General Questions

Gender bias?



Source Bias? (VGG16)

Source Bias: IMDB vs. Wiki



Age Prediction

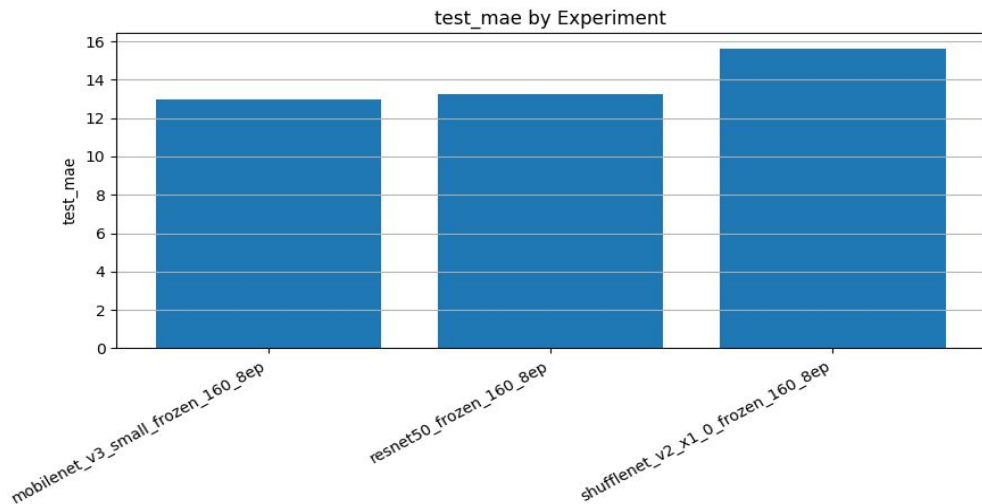
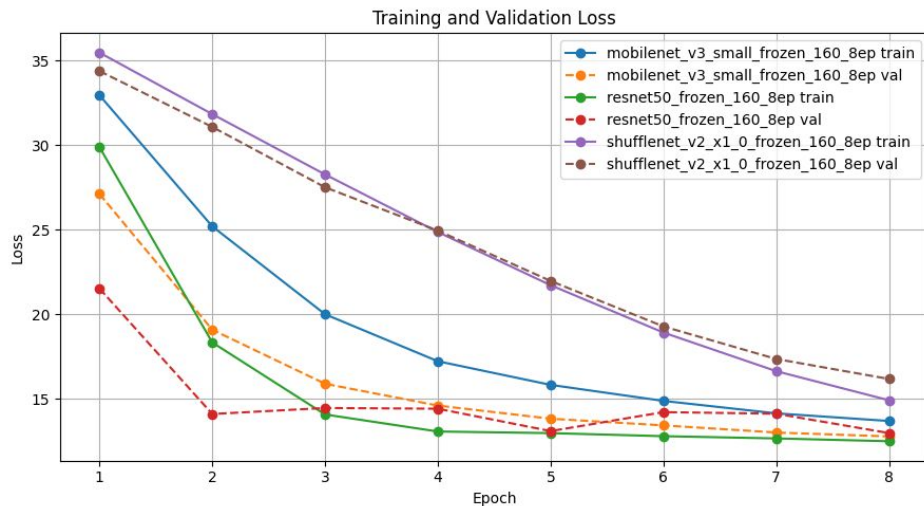
Questions

- Which **pretrained CNN** is best suited?
- How does it do with **only one picture** of each person?
- Training on **genders** separately?

Selection of Model

Model Parameters

- **MobileNet V3 Small**
~2.5 million
- **ShuffleNet V2 x1.0**
~2.3 million
- **ResNet50**
~25.6 million
(4-5 times longer run time)

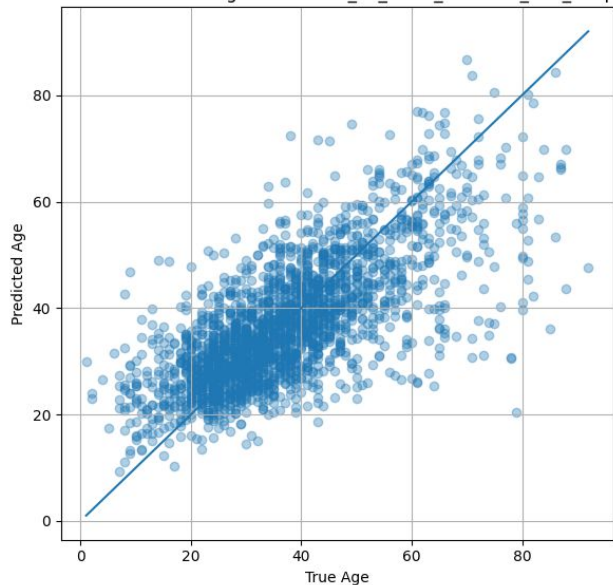


One Picture per Person

Mobilenet

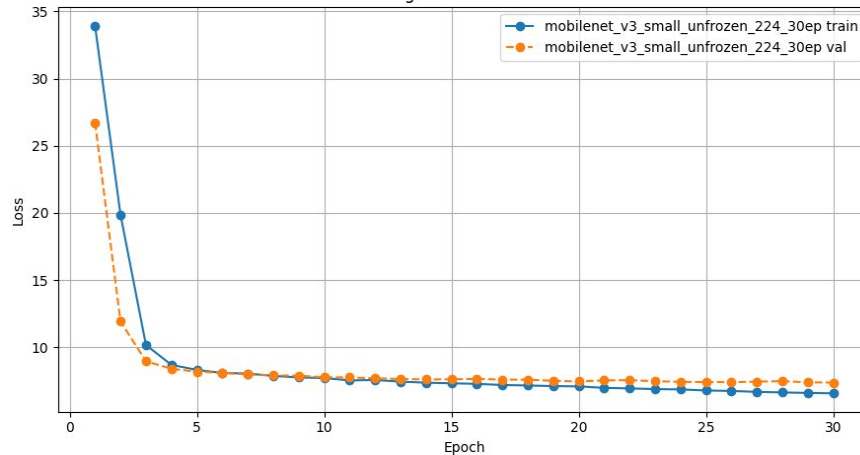
- Successful learning
- Good train-validation agreement
- Little overfitting

Predicted vs True Age: mobilenet_v3_small_unfrozen_224_30ep

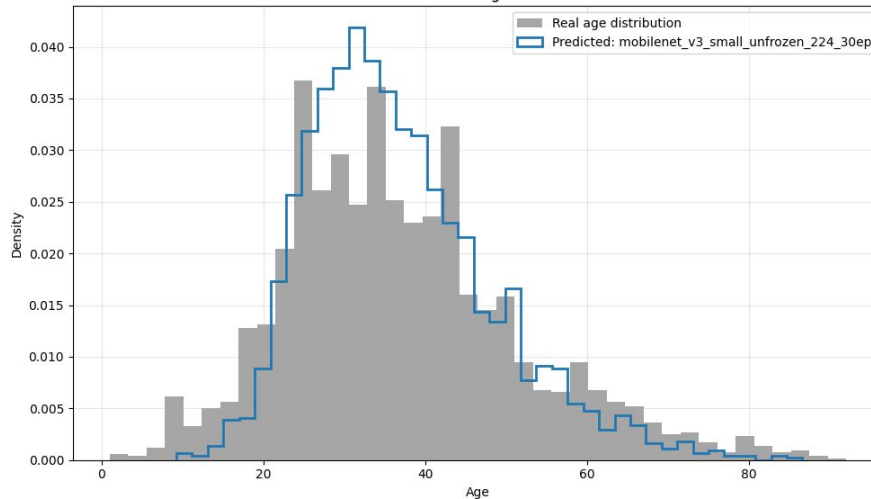


$r^2=0.55$

Training and Validation Loss

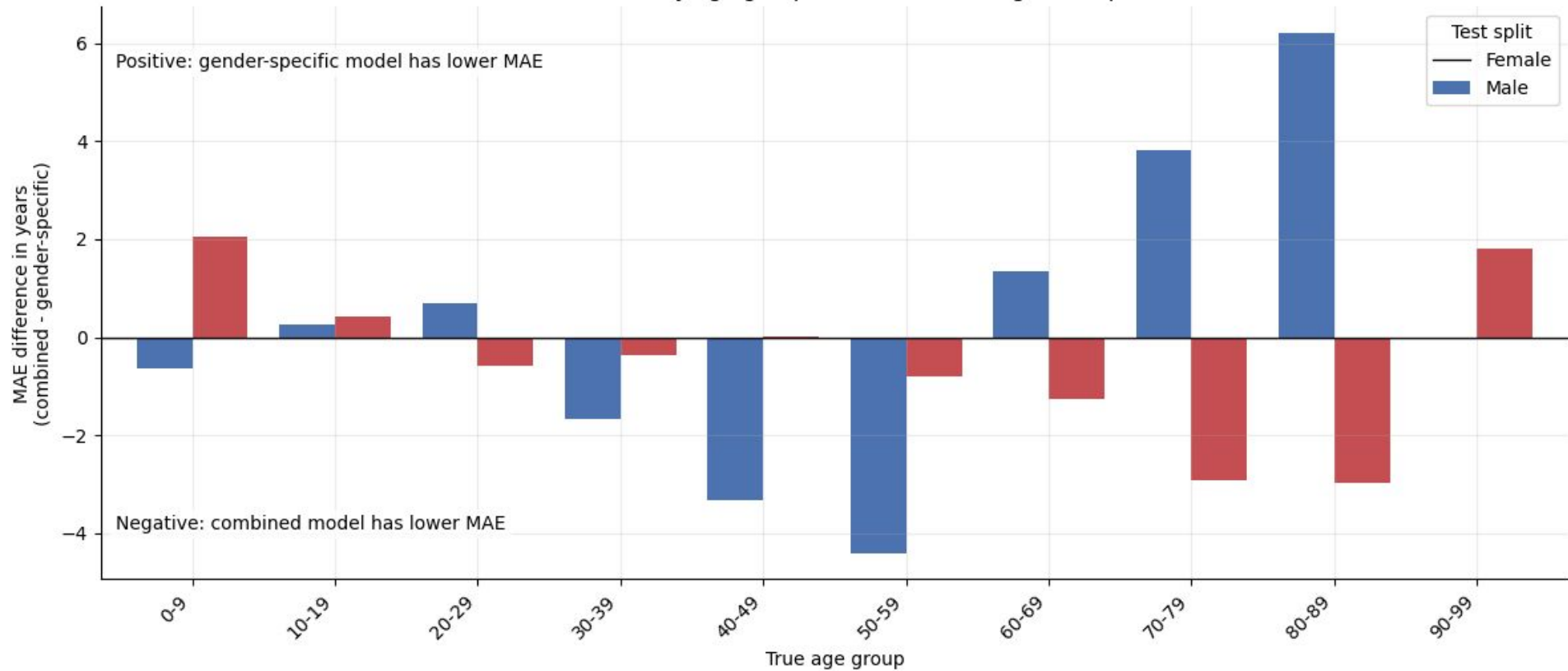


Real vs Predicted Age Distributions



Gender-Specific Training

MAE difference by age group: combined minus gender-specific



Combined model predictions across facial expressions

1Bandit
Predicted age: 31.0



2Cykelhjelm
Predicted age: 48.7



3Solbriller
Predicted age: 47.7



4Hat
Predicted age: 43.4



5Neutral
Predicted age: 39.8



6Glad
Predicted age: 43.2



7Overrasket
Predicted age: 43.7



8Sad
Predicted age: 39.7

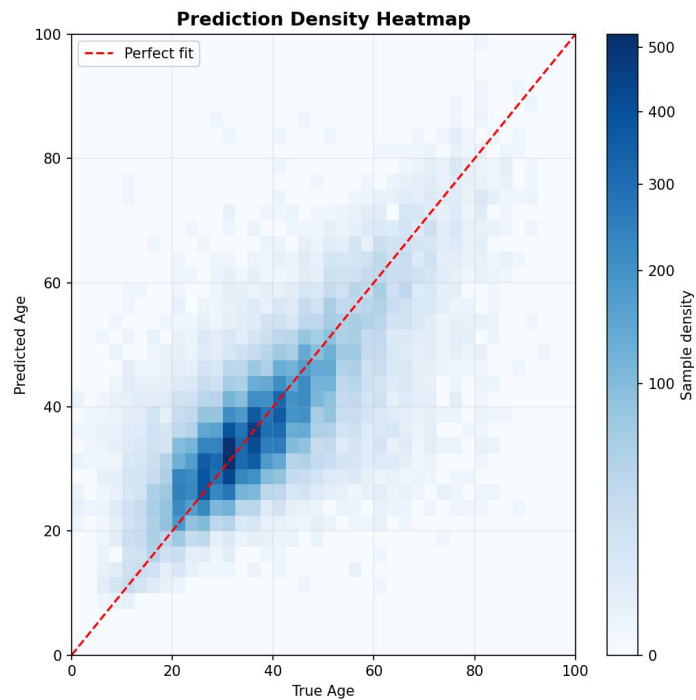


Summary

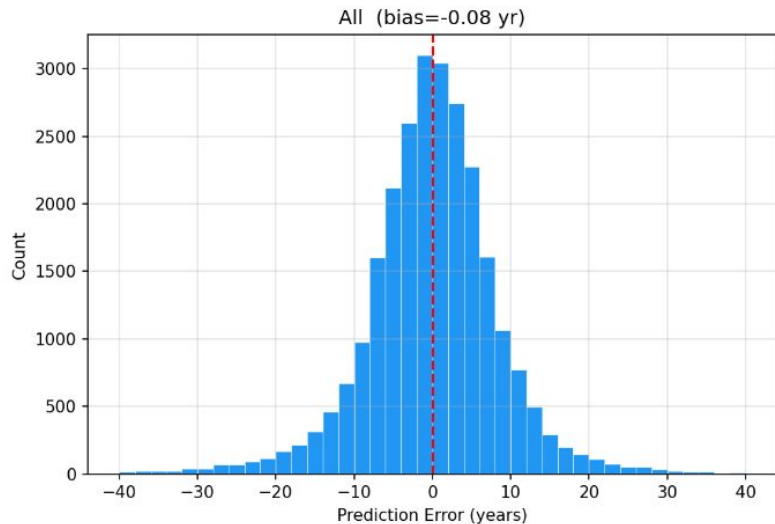
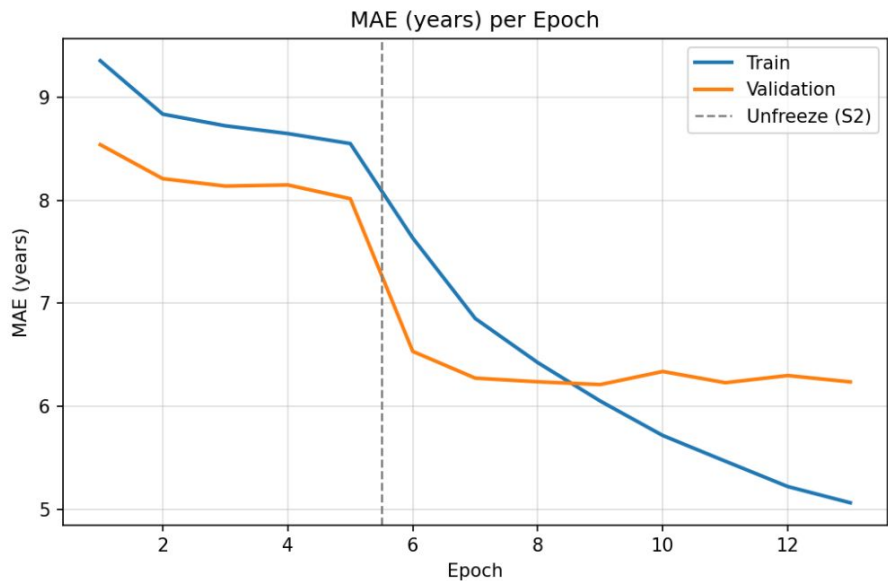
- We decided on Mobilnet and Resnet50
- Age could be predicted with a MAE of 7.3 years
- Training on gender separately didn't help

Age Regression Using ResNet50

- Regression head
- **Stage 1:** backbone fully frozen, head-only training (~ 1M Parameters)
- **Stage 2:** unfreeze Layer 3 + Layer 4 at once (~ 23M params)



Age Regression Using ResNet50



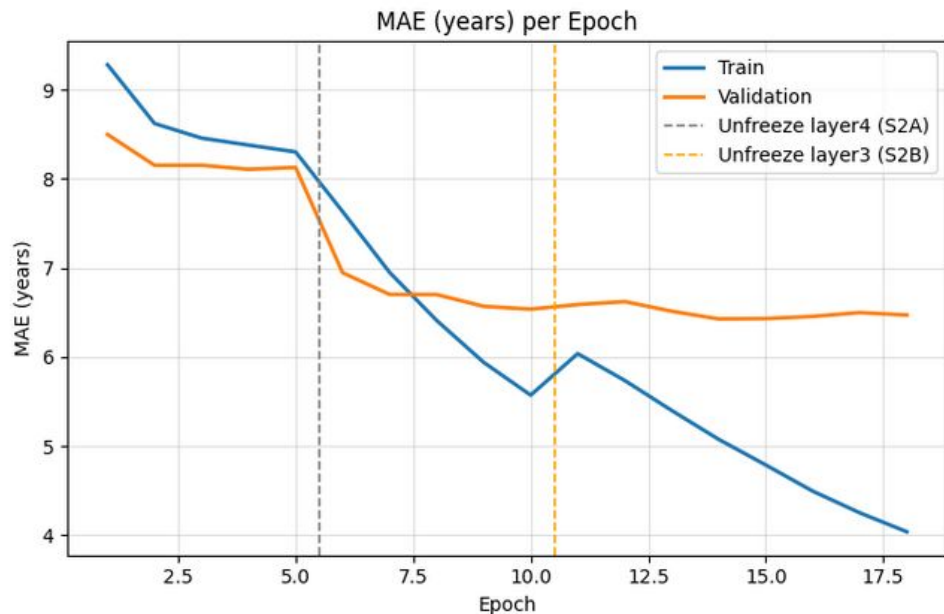
Attempts at Improving Previous Model

- **Stage 1:** same
- **Stage 2A:** unfreeze Layer 4 + head only (~6M params)
- **Stage 2B:** unfreeze Layer 3 + Layer 4 + head (~23M params), same differential LR

⇒ ResNet50 needs minimal training to adapt

⇒ Overfitting bad labels

⇒ Limit seems to be around MAE ~ 6 yrs



Best and Worst Age Predictions

ResNet50 - 5 Best Predictions

True: 27
Pred: 27.0
|err|=0.0



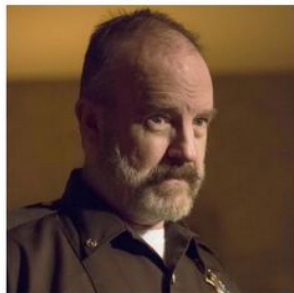
True: 30
Pred: 30.0
|err|=0.0



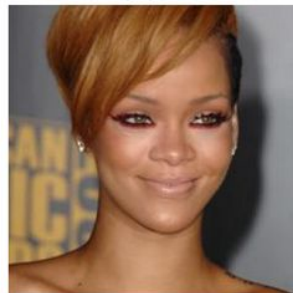
True: 35
Pred: 35.0
|err|=0.0



True: 60
Pred: 60.0
|err|=0.0



True: 21
Pred: 21.0
|err|=0.0



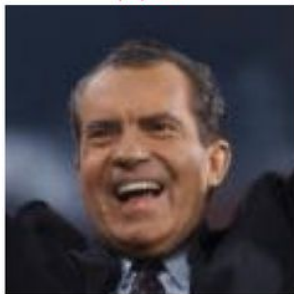
Saved: ResNet50_best_predictions.png

ResNet50 - 5 Worst Predictions

True: 89
Pred: 21.0
|err|=68.0



True: 94
Pred: 35.2
|err|=58.8



True: 26
Pred: 26.3
|err|=57.7



True: 79
Pred: 21.6
|err|=57.4



True: 12
Pred: 68.8
|err|=56.8



Saved: ResNet50_worst_predictions.png

Classification

Age Prediction

Age Bin Classification Examples

- Same backbone and staged fine-tuning
- Head outputs 20 logits (age-bins)
- Predicted age decoded as expected value

True: 15 yr
Pred: 25.2 yr



True: 25 yr
Pred: 34.9 yr



True: 35 yr
Pred: 32.8 yr



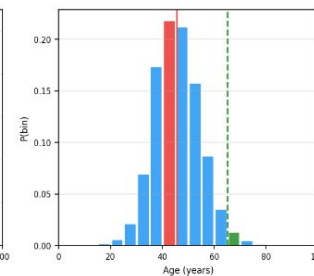
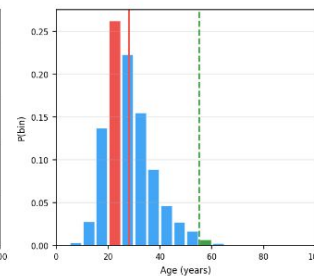
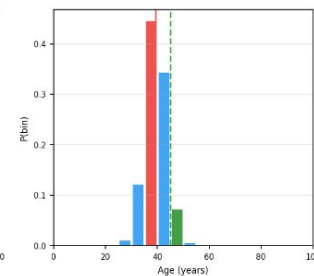
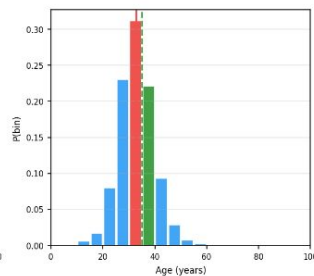
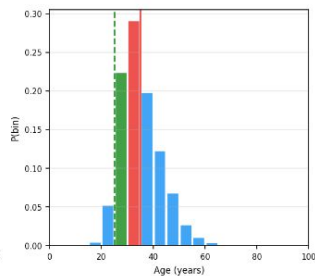
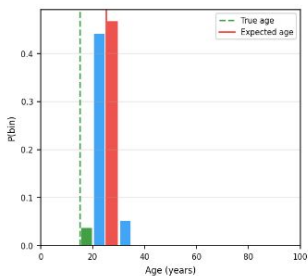
True: 45 yr
Pred: 39.3 yr



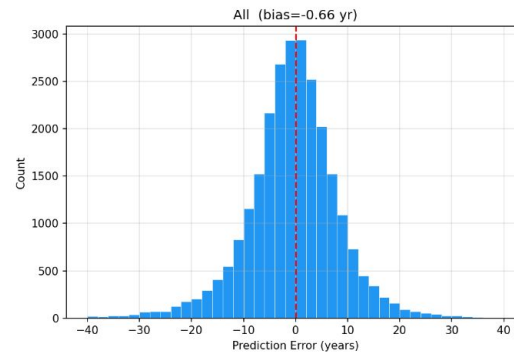
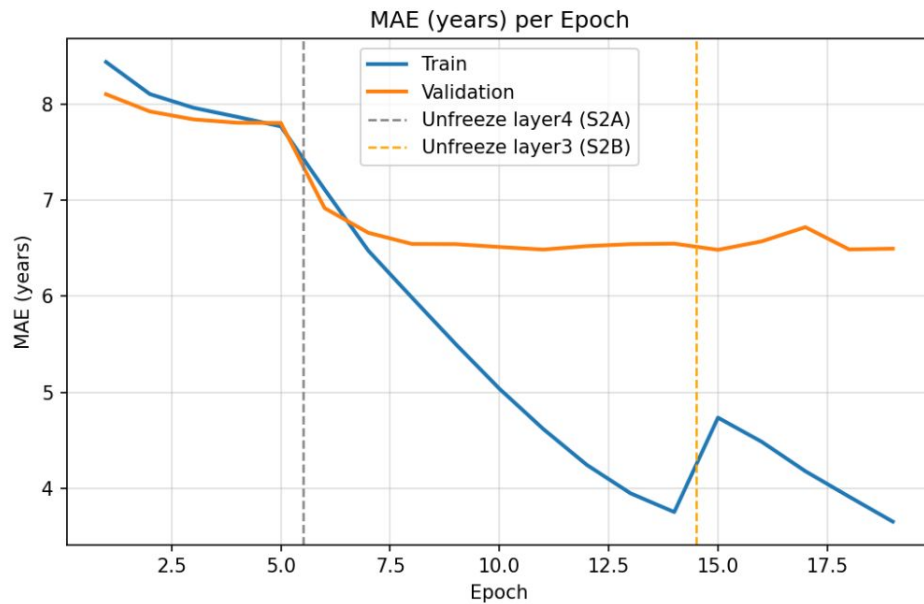
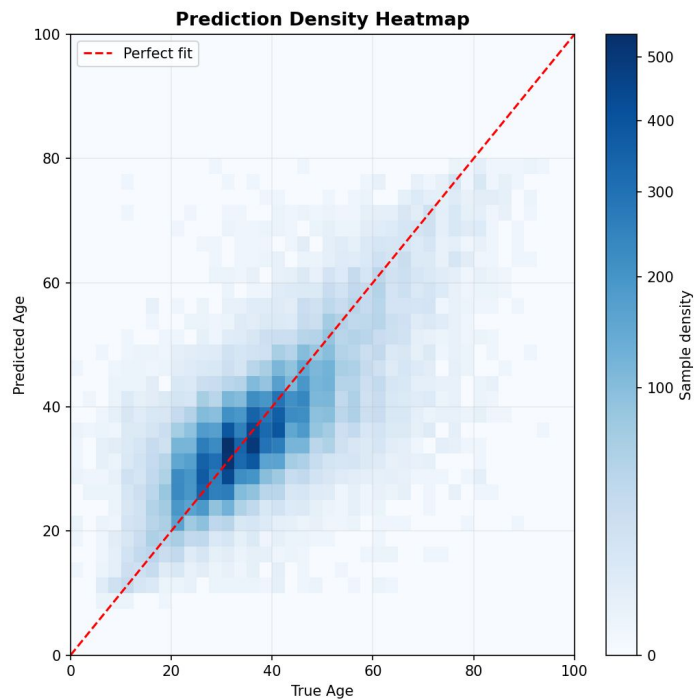
True: 55 yr
Pred: 28.0 yr



True: 65 yr
Pred: 45.6 yr



Age-Bin Classification



Gender Prediction

Gender Classification

3 Models:

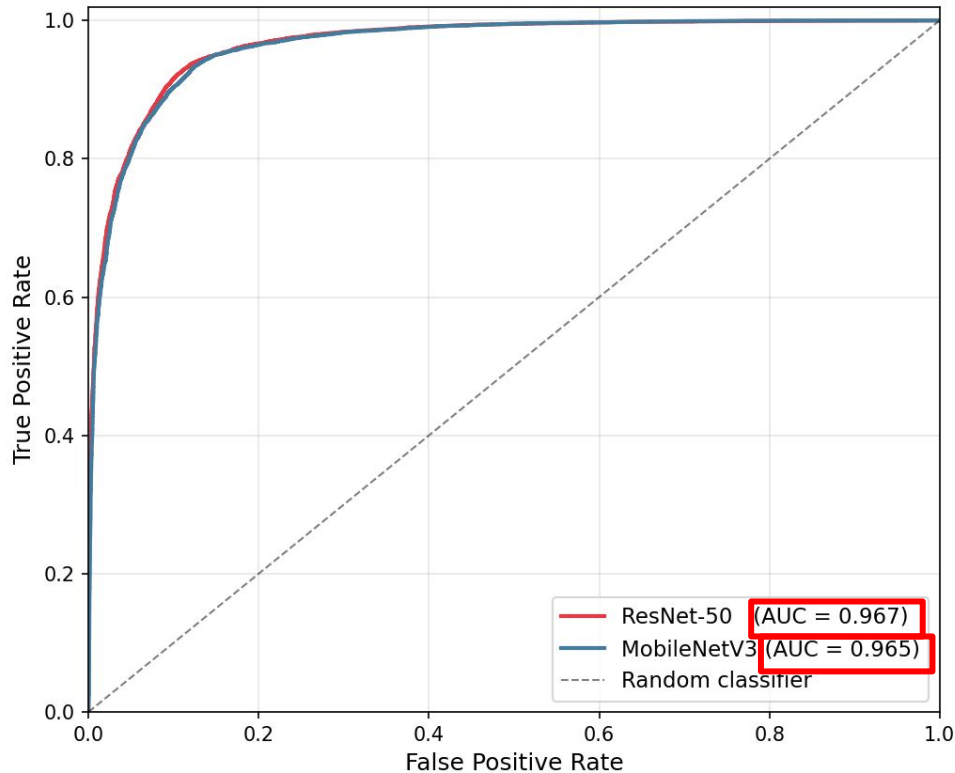
- ResNet50
- MobileNet
- non pretrained CNN

Goal:

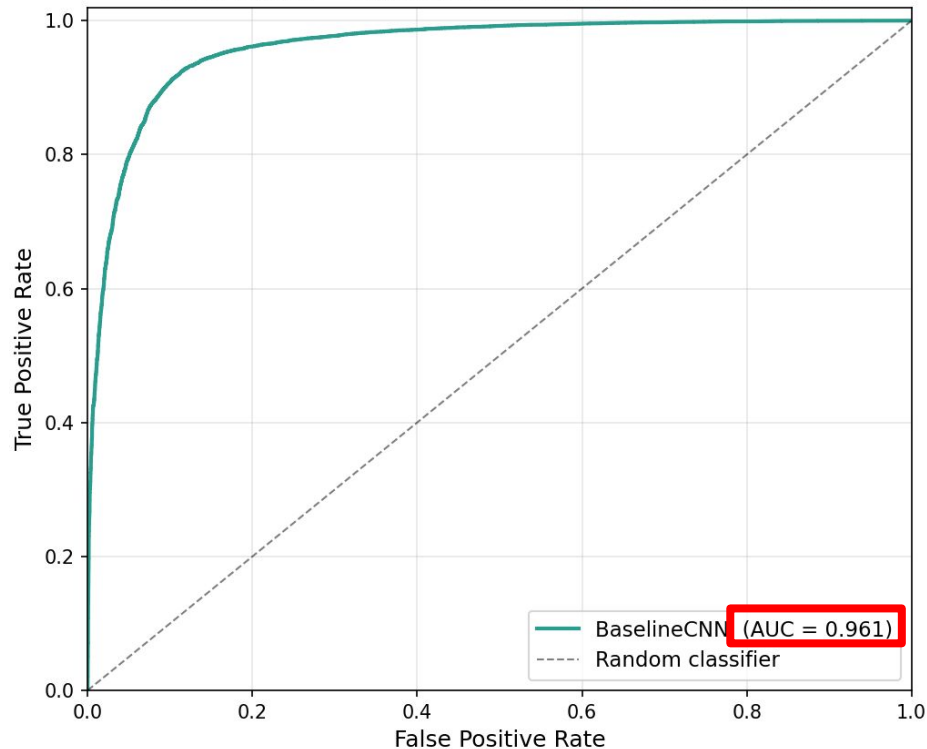
- Predict gender based on labeled data
- Study limitations of the models and possible improvements

Gender Classification - Model Evaluation

ROC Curve — Gender Classification

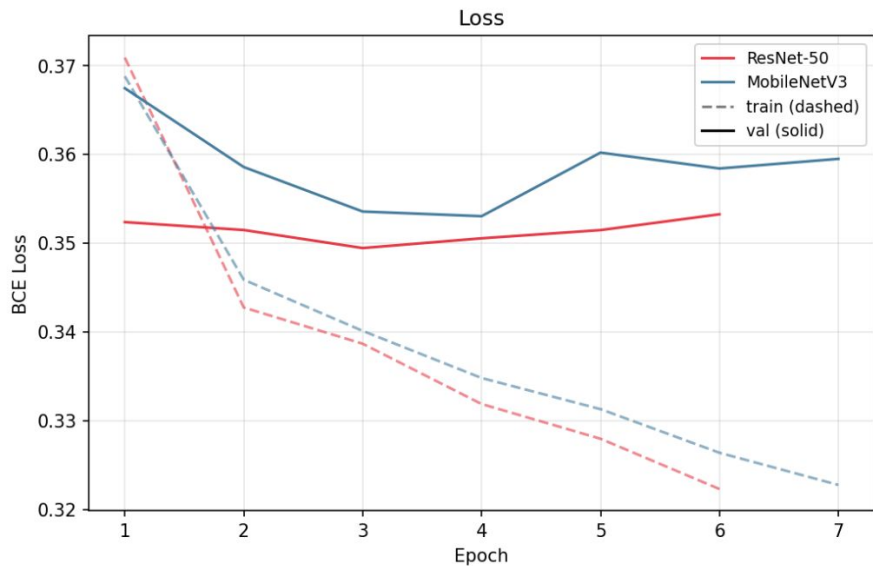


ROC Curve — BaselineCNN Gender Classification

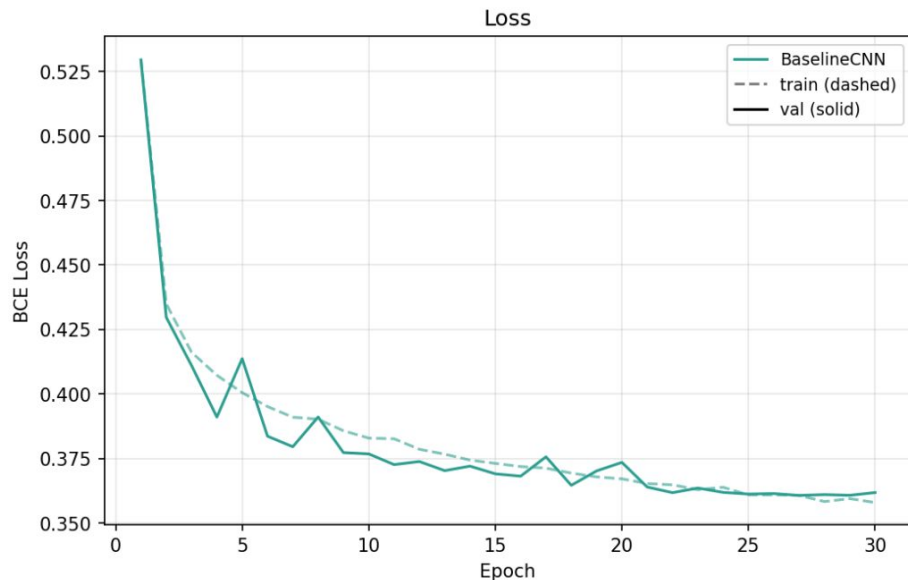


Gender Classification - Model Evaluation

ResNet-50 and MobileNetV3



Baseline CNN



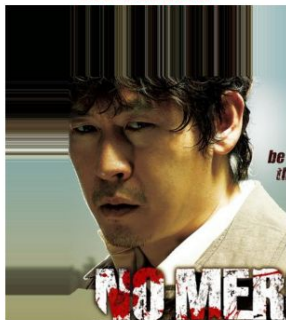
Gender Classification - Model Evaluation

correctly classified:

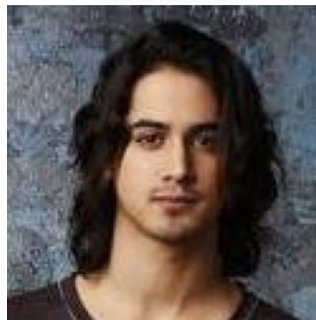
Female 94% GT: Female



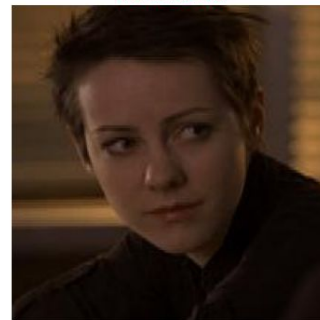
True: Male
Pred: Female (93.3%)



True: Male
Pred: Female (92.4%)



True: Female
Pred: Male (94.4%)



True: Female
Pred: Male (96.6%)



correctly classified:

Male 95% GT: Male



Grad-CAM = Gradient-weighted Class Activation Mapping

Discussion

Summary

Clustering:

- Investigated patterns

Age Prediction:

- Best prediction with MAE 6 years

Gender Prediction:

- Best prediction with AUC 0.967

Discussion

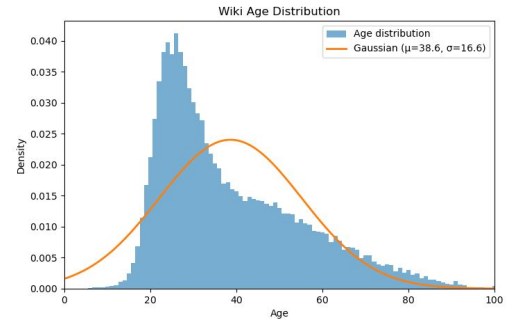
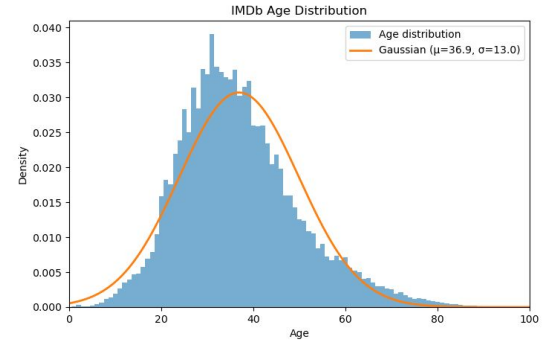
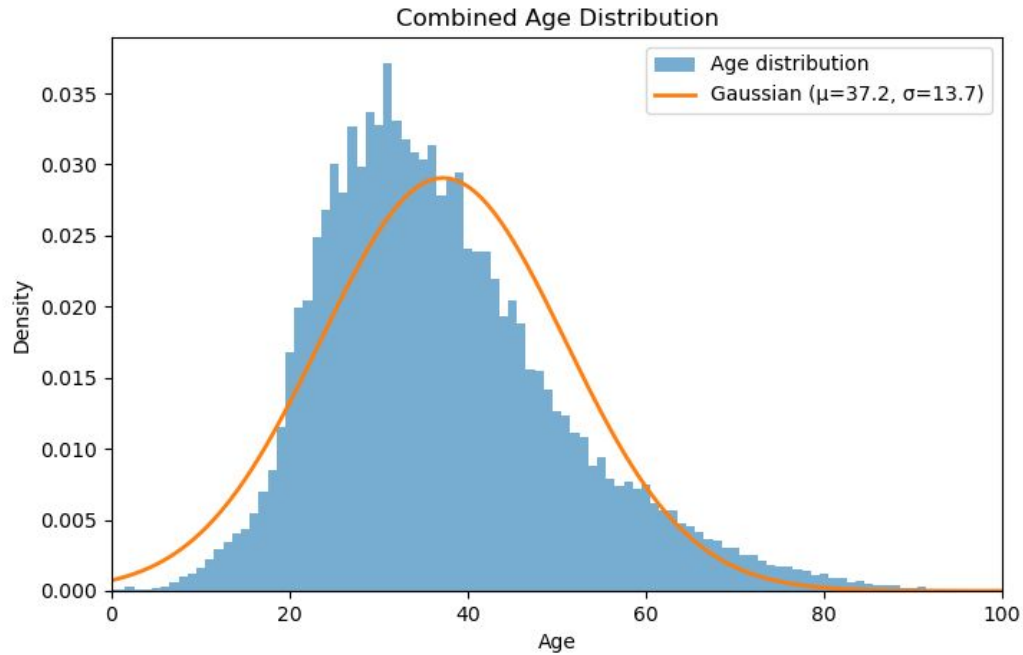
Our Takeaways:

- Load your images effectively!
- Use a GPU! Shoutout to Kaggle Notebooks!
- Watch for noisy labels and data bias!
- Regression and Classification Models:
 - ResNet50 > MobileNet V3 > ShuffleNet V2 > non pretrained CNN

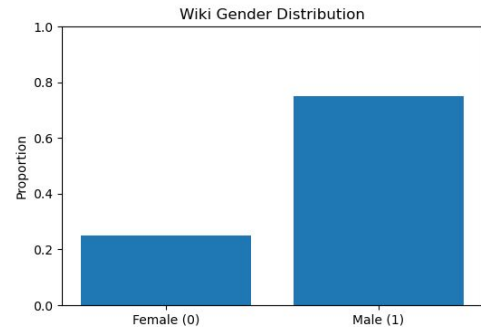
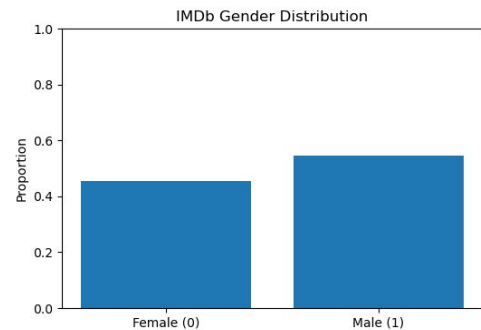
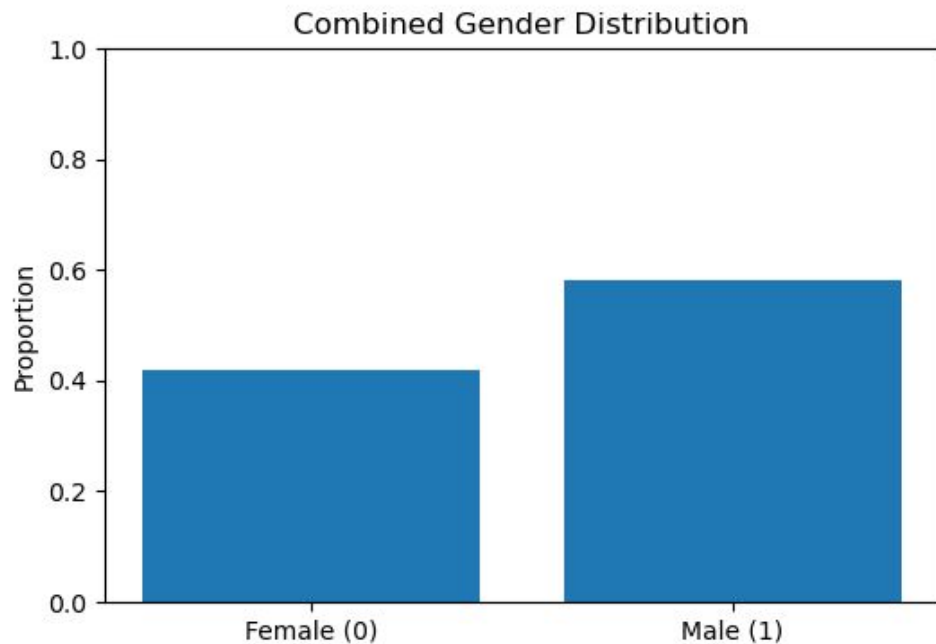
Appendix

Appendix - Data

The IMDB-WIKI Dataset - Age Distribution



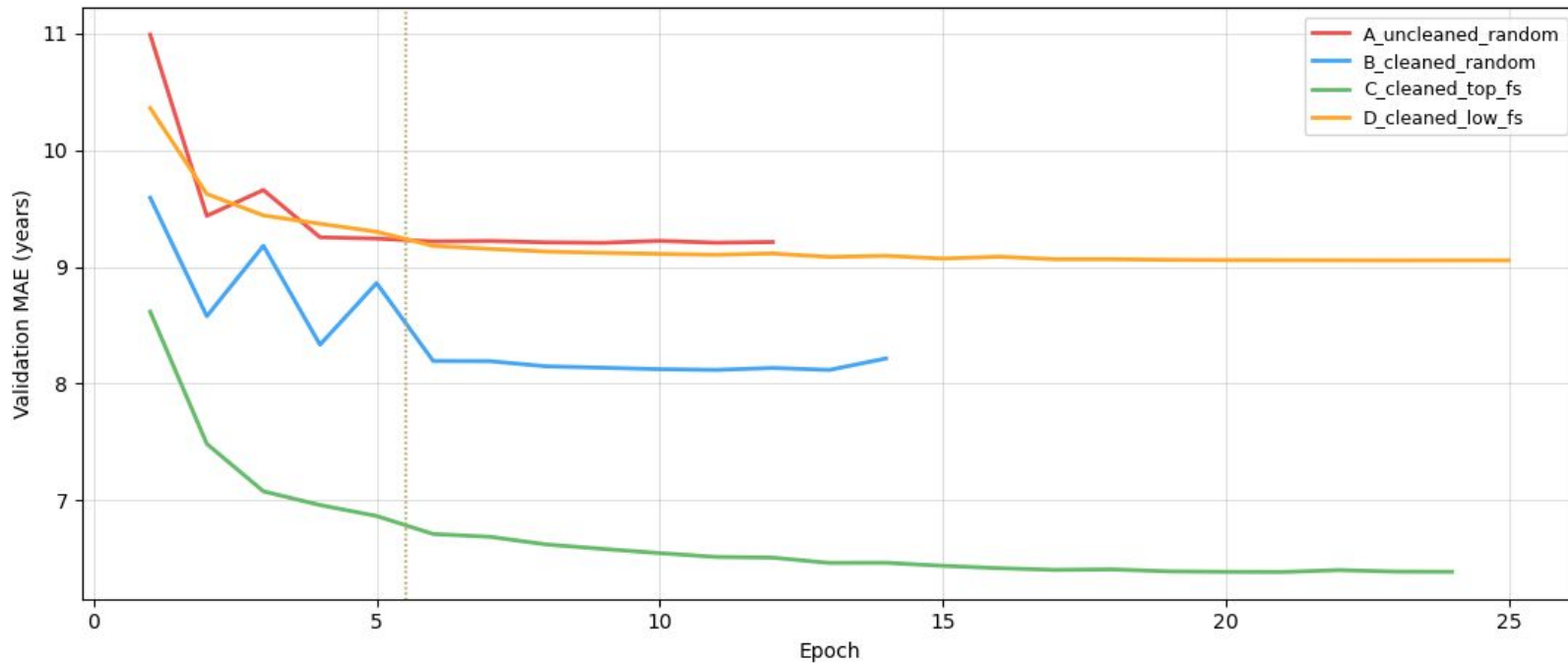
The IMDB-WIKI Dataset - Gender Distribution



Experiment on Data Quality

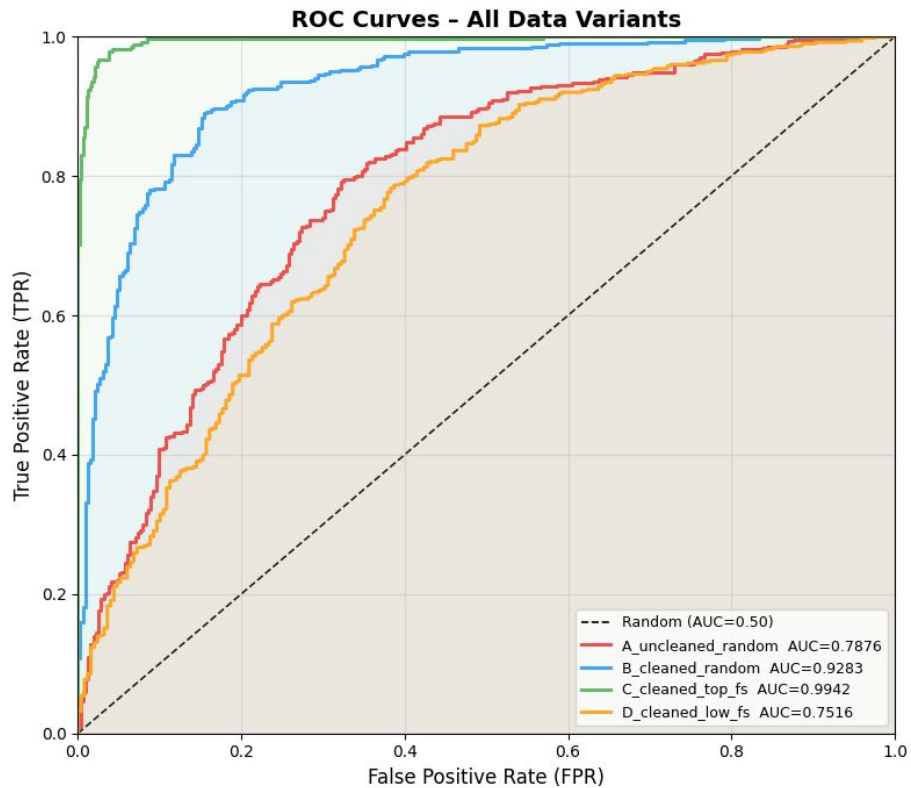
light mobilenet
model

Validation MAE - All Variants



Experiment on Data Quality

light mobilenet
model



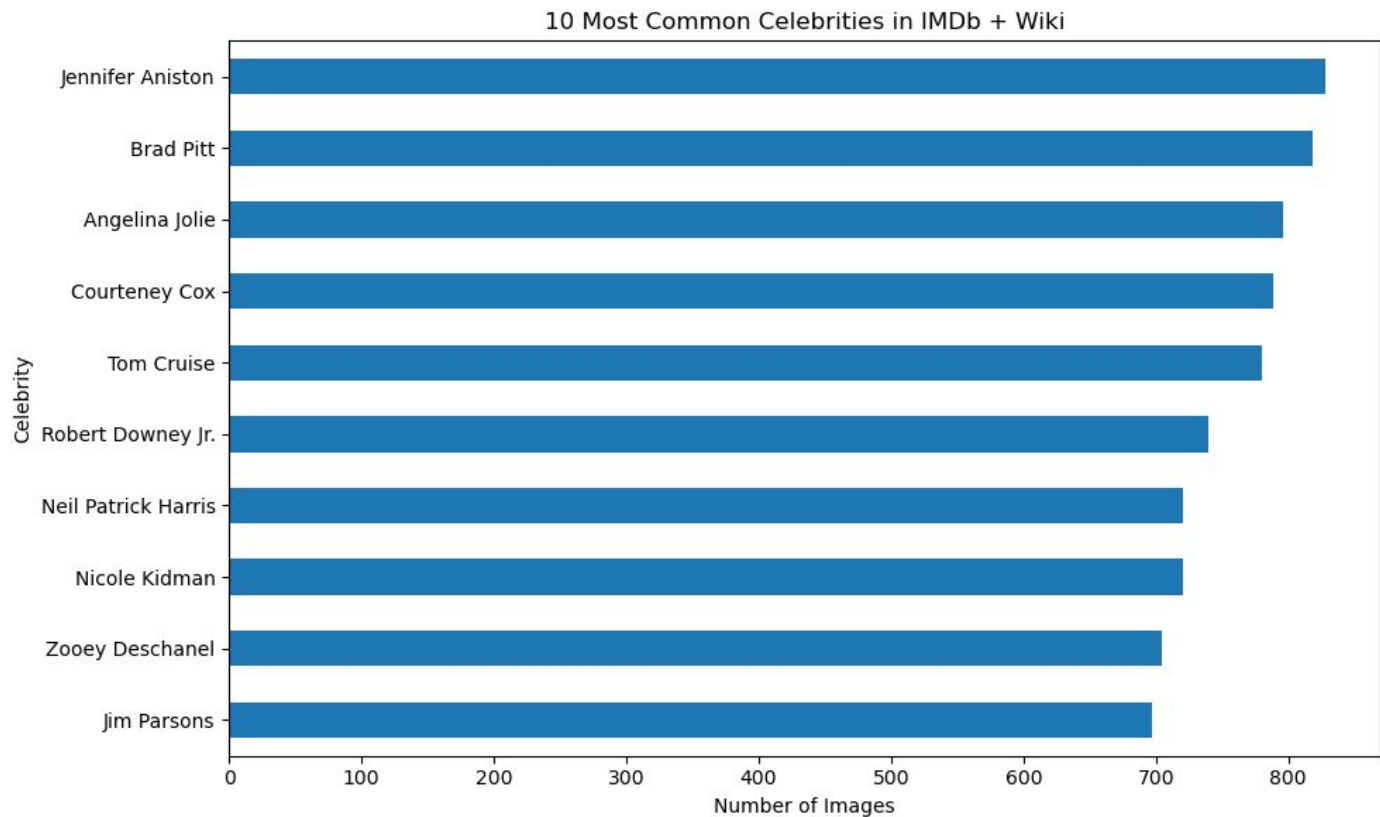
The IMDB-WIKI Dataset - Overview

- Large dataset of celebrities
- Crawled the IMDb website for all images, date of birth & gender
- Crawled all profile images from pages of people from Wikipedia with the same meta information
- Removed all images without a timestamp
- Images can be from movies with extended production time, wrong time stamps
- IMDb 460,723 face images from 20,284 celebrities
- Wikipedia 62,328

The IMDB-WIKI Dataset - Meta Information

- **dob:** date of birth (Matlab serial date number)
- **photo_taken:** year when the photo was taken
- **full_path:** path to file
- **gender:** 0 for female and 1 for male, NaN if unknown
- **name:** name of the celebrity
- **face_location:** location of the face. To crop the face in Matlab run
- **face_score:** detector score (the higher the better). *Inf* implies that no face was found in the image and the *face_location* then just returns the entire image
- **second_face_score:** detector score of the face with the second highest score. This is useful to ignore images with more than one face. *second_face_score* is *NaN* if no second face was detected
- **celeb_names:** list of all celebrity names
- **celeb_id:** index of celebrity name

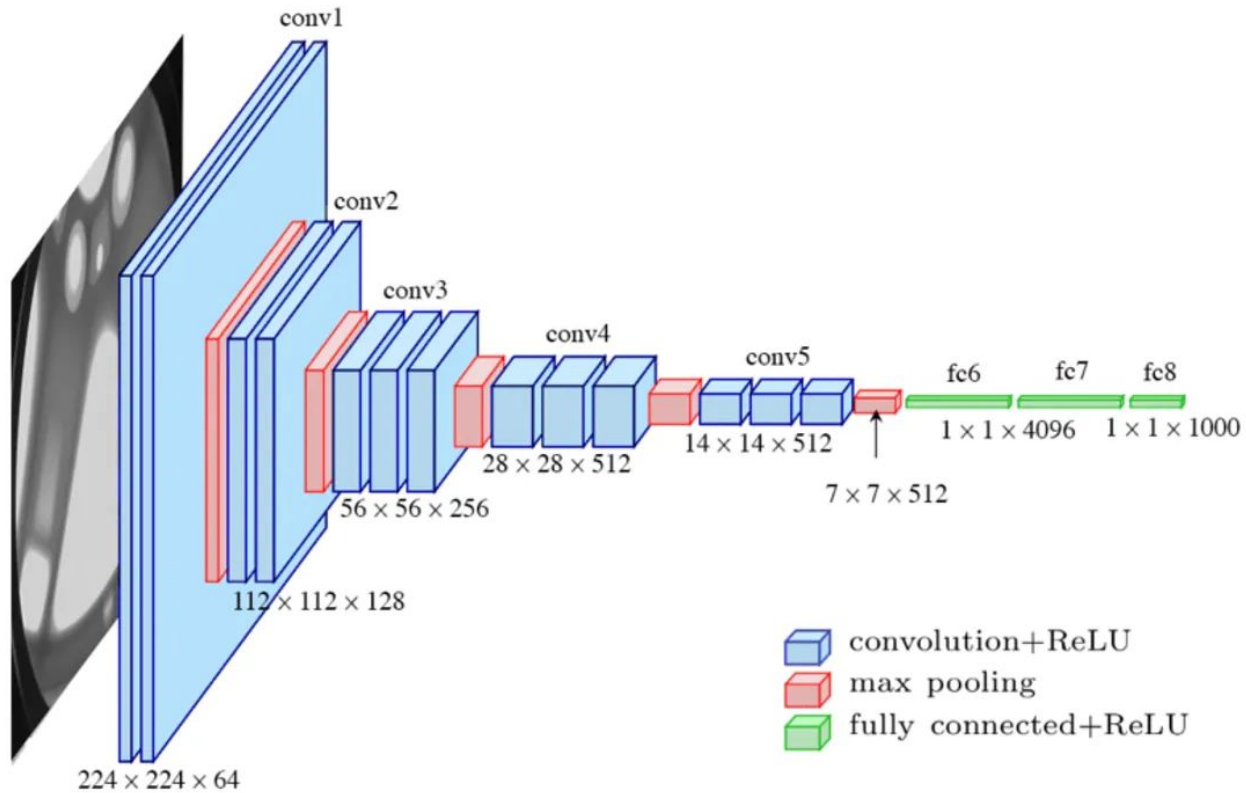
The IMDB-WIKI Dataset - Multiplicity



Appendix - Clustering

VGG 16

VGG16 - Architecture



Looking at the Clusters (Clear: People in Suits)

Cluster 2 - \bar{x} Age 50yrs ($\sigma=19.1$ yrs)

16yrs



28yrs



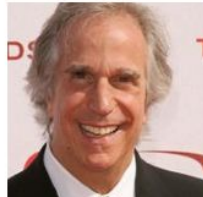
14yrs



64yrs



63yrs



61yrs



50yrs



52yrs



19yrs



48yrs



30yrs



47yrs



78yrs



60yrs



55yrs



36yrs



28yrs



53yrs



70yrs



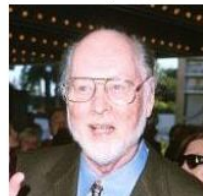
57yrs



71yrs



68yrs



55yrs



78yrs



Looking at the Clusters (Clear: Sports People)

Cluster 5 - \emptyset Age 27yrs ($\sigma=10.1$ yrs)

31yrs



46yrs



23yrs



35yrs



24yrs



22yrs



20yrs



19yrs



20yrs



17yrs



20yrs



36yrs



16yrs



44yrs



18yrs



54yrs



19yrs



28yrs



24yrs



27yrs



33yrs



22yrs



22yrs

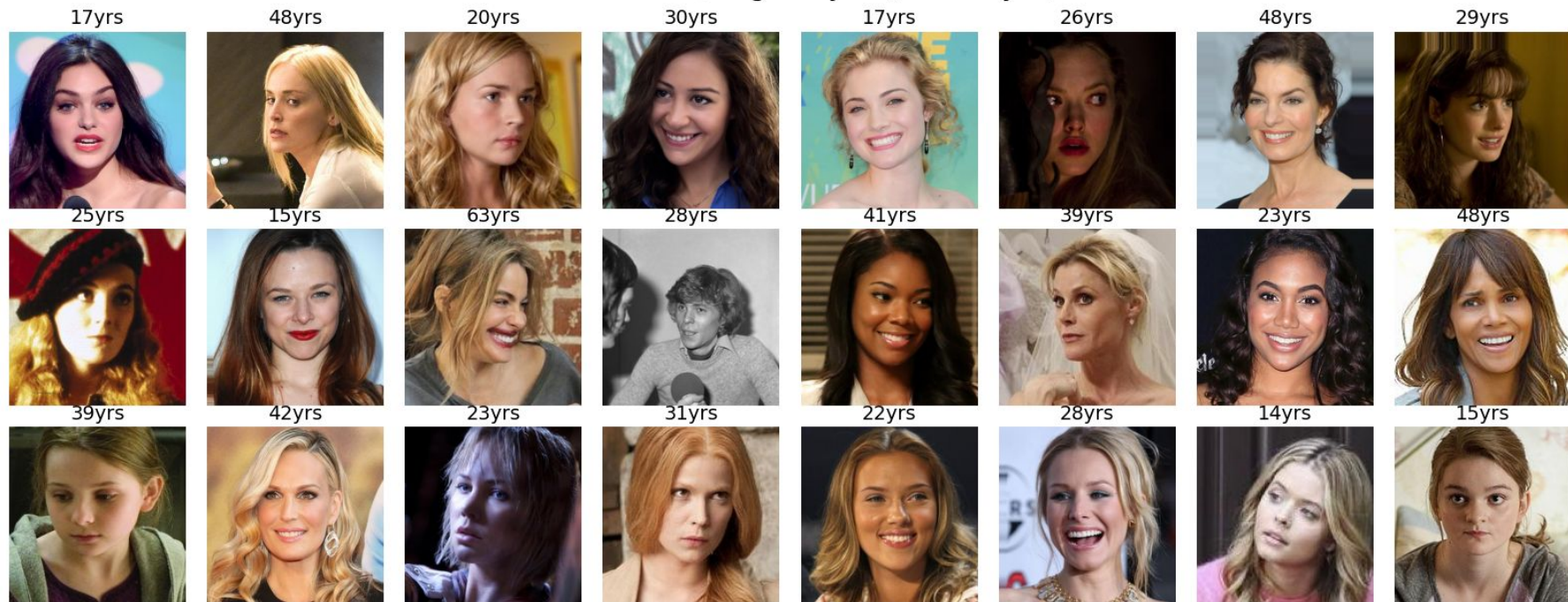


19yrs



Looking at the (not so clear) Clusters

Cluster 1 - $\bar{\text{Age}}$ 30yrs ($\sigma=12.9\text{yrs}$)



Looking at the (not so clear) Clusters

Cluster 4 - \bar{X} Age 30yrs ($\sigma=15.0$ yrs)



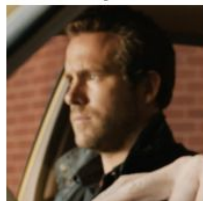
Looking at the (not so clear) Clusters

Cluster 6 - \bar{X} Age 47yrs ($\sigma=20.8$ yrs)

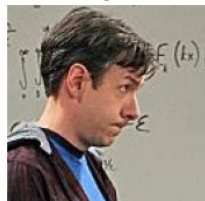
24yrs



64yrs



34yrs



61yrs



60yrs



87yrs



17yrs



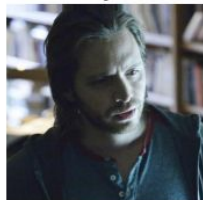
75yrs



63yrs



37yrs



29yrs



44yrs



19yrs



54yrs



48yrs



11yrs



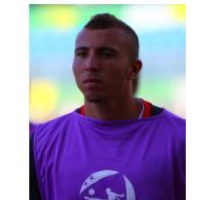
66yrs



63yrs



18yrs



42yrs



46yrs



79yrs



49yrs



45yrs



Looking at the (not so clear) Clusters

Cluster 3 - $\bar{\text{Age}}$ 43yrs ($\sigma=20.5\text{yrs}$)

38yrs



53yrs



15yrs



81yrs



26yrs



49yrs



53yrs



57yrs



59yrs



25yrs



43yrs



16yrs



92yrs



11yrs



50yrs



24yrs



52yrs



31yrs



31yrs



27yrs



52yrs



40yrs



72yrs



43yrs



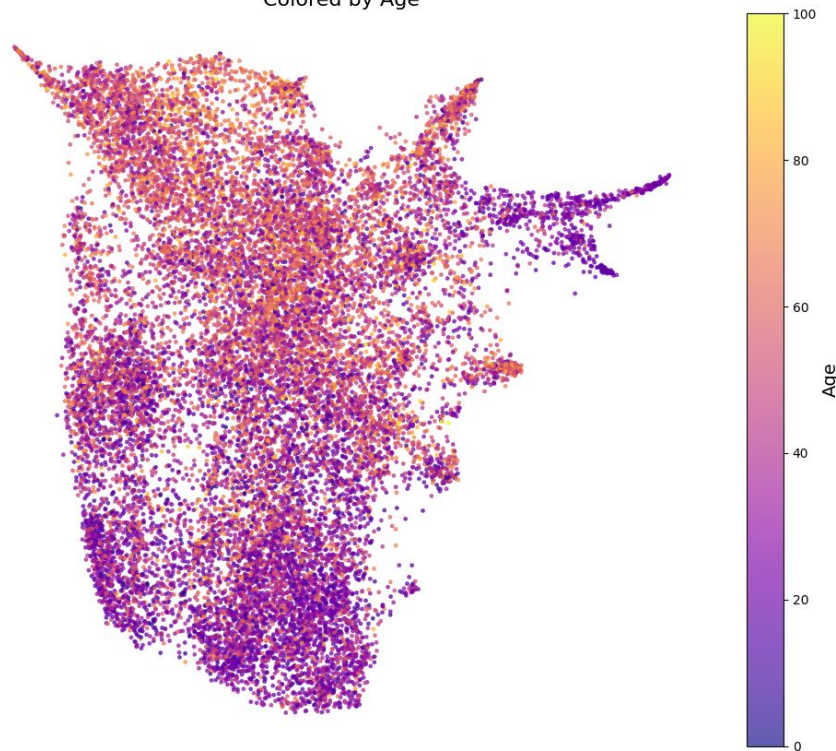
VGG16 (K-Means) - Visualization by UMAP

VGG16 Features - UMAP

Cluster & Thumbnails

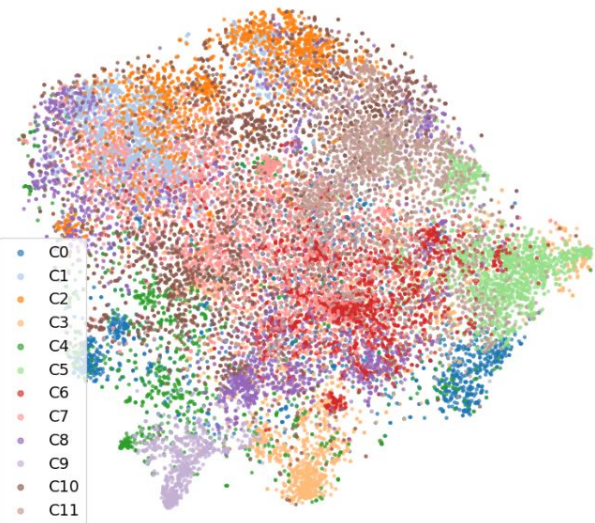


Colored by Age



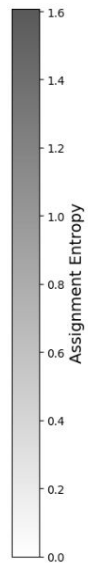
GMM Clustering on VGG16

GMM Clusters

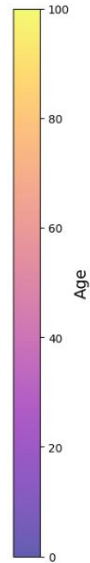


GMM (k=12) - t-SNE Visualization

GMM Uncertainty
(high = point sits between clusters)

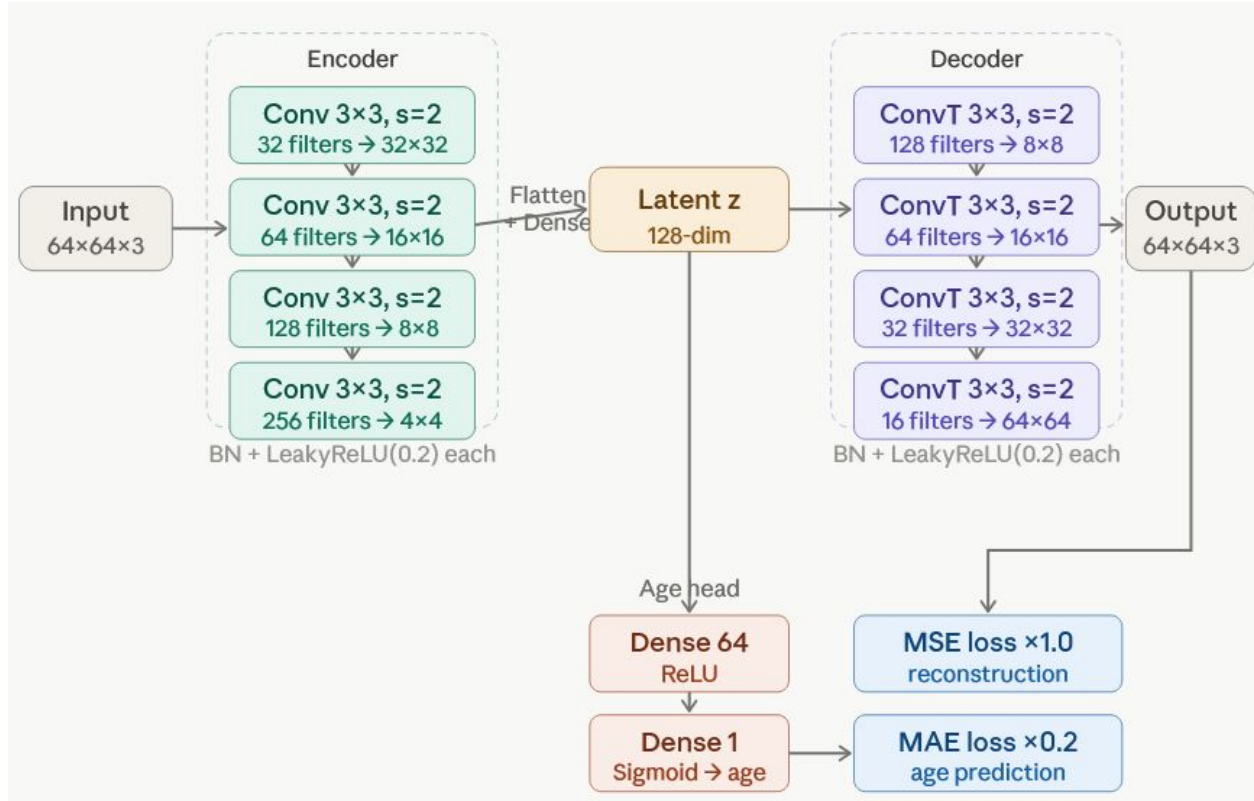


Age (for reference)



Age-Supervised Autoencoder

Age-Supervised AE (SAE) - Architecture



Age-Supervised AE (SAE) - Model Description

Model ready | AGE_LOSS_WEIGHT=0.2

Model: "age_supervised_ae"

Layer (type)	Output Shape	Param #	Connected to
input_layer_2 (InputLayer)	(None, 64, 64, 3)	0	-
encoder (Functional)	(None, 128)	914,752	input_layer_2[0][0]
dense_1 (Dense)	(None, 64)	8,256	encoder[0][0]
age_pred (Dense)	(None, 1)	65	dense_1[0][0]
decoder (Functional)	(None, 64, 64, 3)	921,699	encoder[0][0]

Total params: 1,844,772 (7.04 MB)

Trainable params: 1,843,332 (7.03 MB)

Non-trainable params: 1,440 (5.62 KB)

Looking at the Clusters (SAE works well)

Cluster 6 - \bar{x} Age 66yrs ($\sigma=10.9$ yrs)

55yrs



48yrs



64yrs



76yrs



55yrs



63yrs



63yrs



57yrs



74yrs



62yrs



90yrs



61yrs



70yrs



50yrs



70yrs



74yrs



84yrs



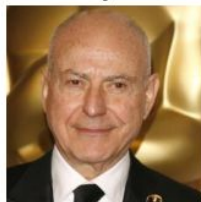
66yrs



49yrs



73yrs



73yrs



66yrs



71yrs



80yrs



Looking at the Clusters (SAE works not so well)

Cluster 2 - \bar{x} Age 32yrs ($\sigma=16.8$ yrs)

20yrs



23yrs



19yrs



26yrs



20yrs



36yrs



21yrs



47yrs



7yrs



56yrs



51yrs



20yrs



23yrs



44yrs



16yrs



64yrs



20yrs



20yrs



19yrs



33yrs



63yrs



33yrs



21yrs

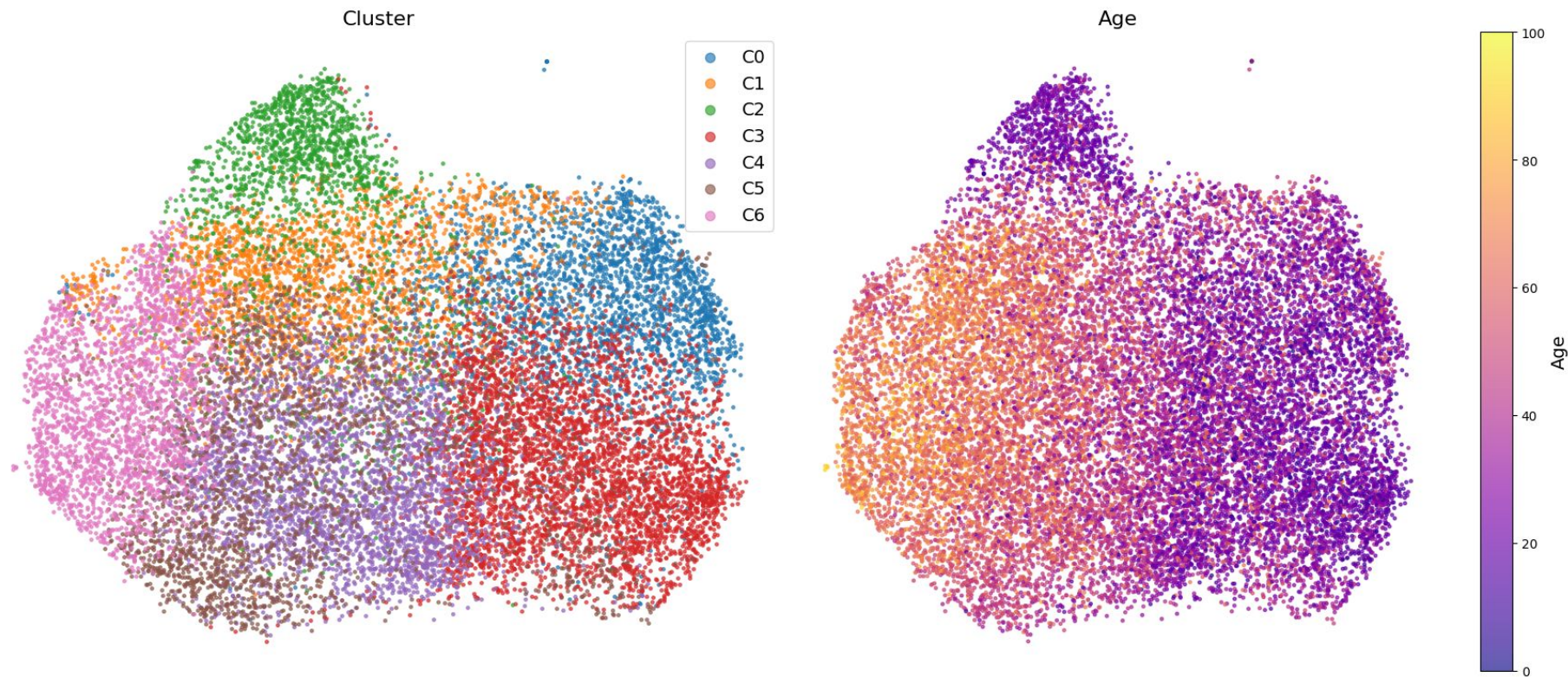


61yrs



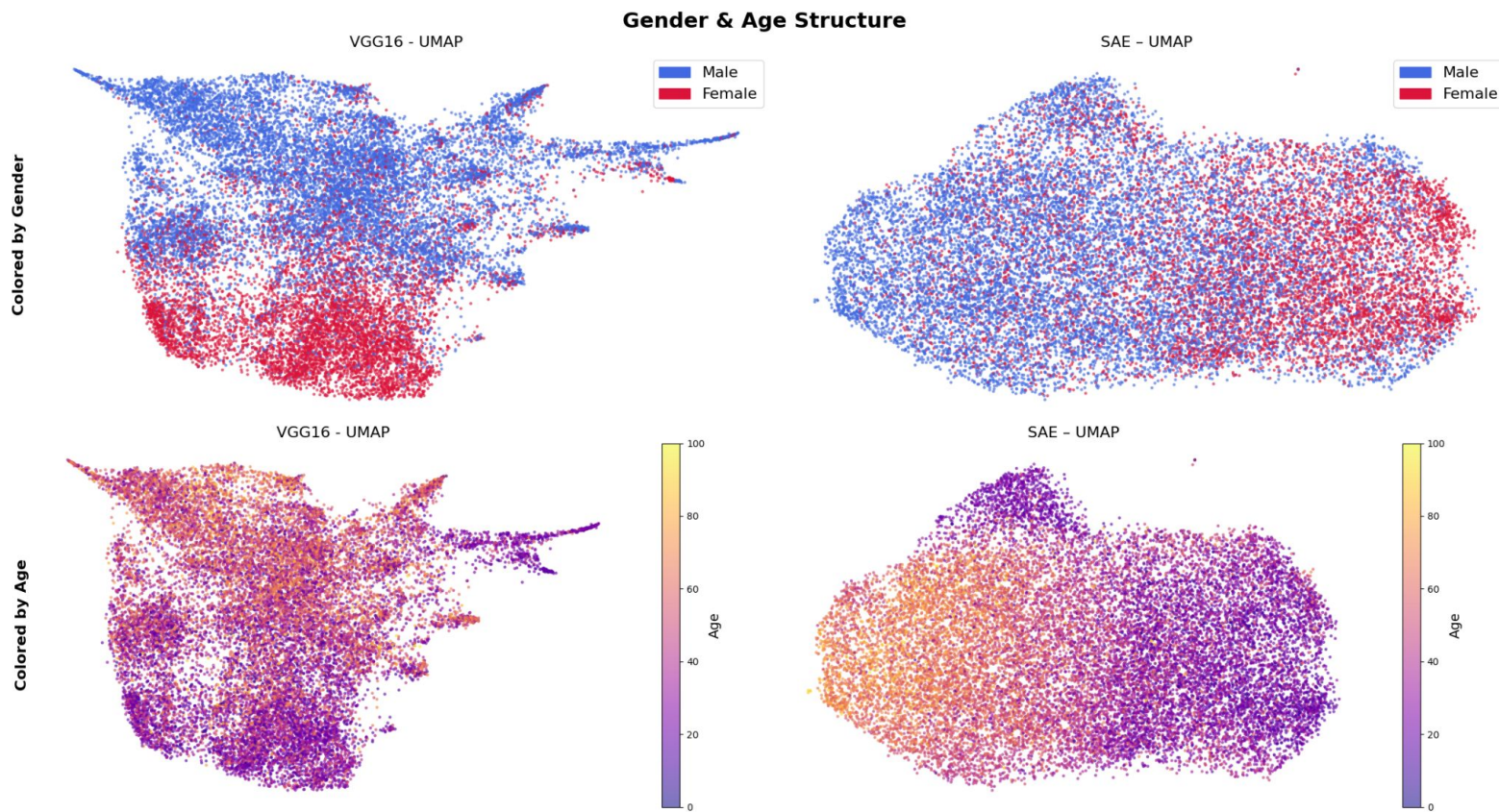
SAE - Visualization by UMAP

Age-Supervised AE - UMAP



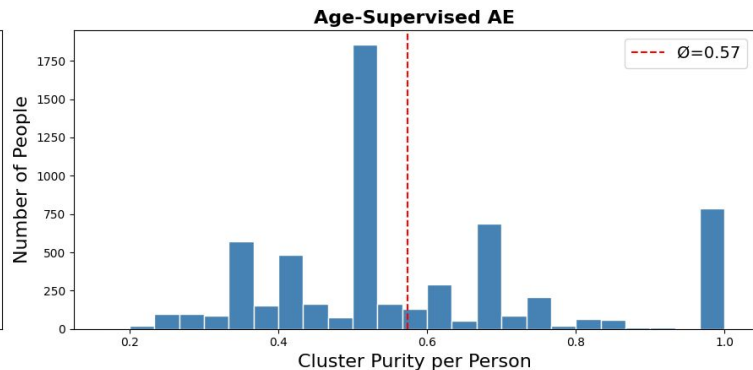
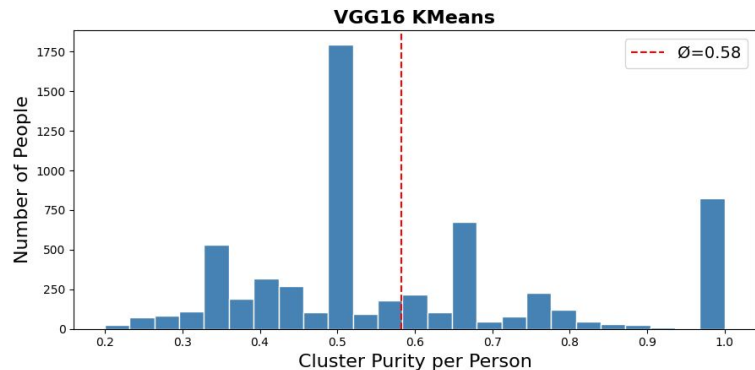
General Questions

Gender bias? - UMAP Visualization

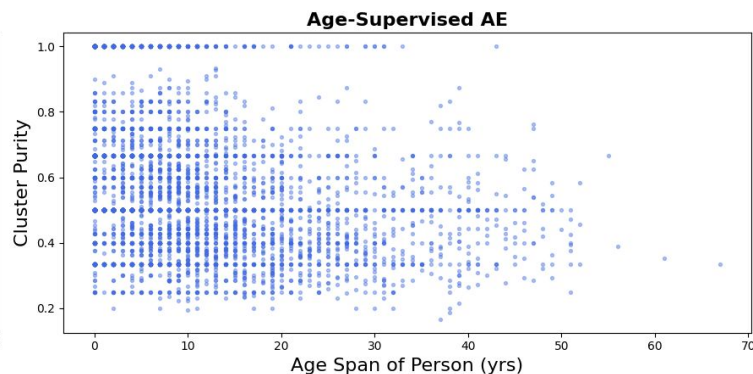
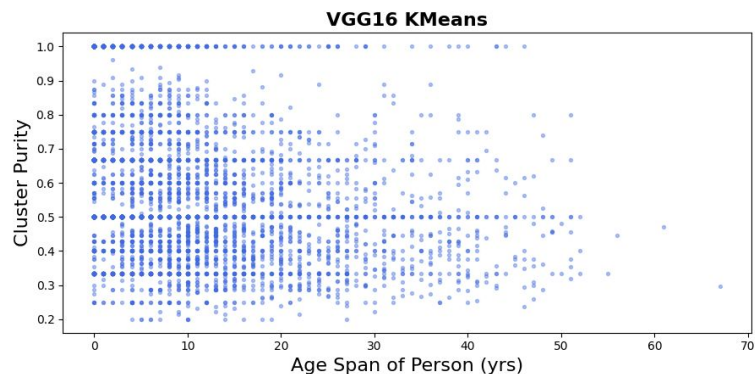


Images of the Same Person?

Person Consistency: Do same-person images land in the same cluster?



Purity vs. Age Span: Does appearance change across age hurt consistency?



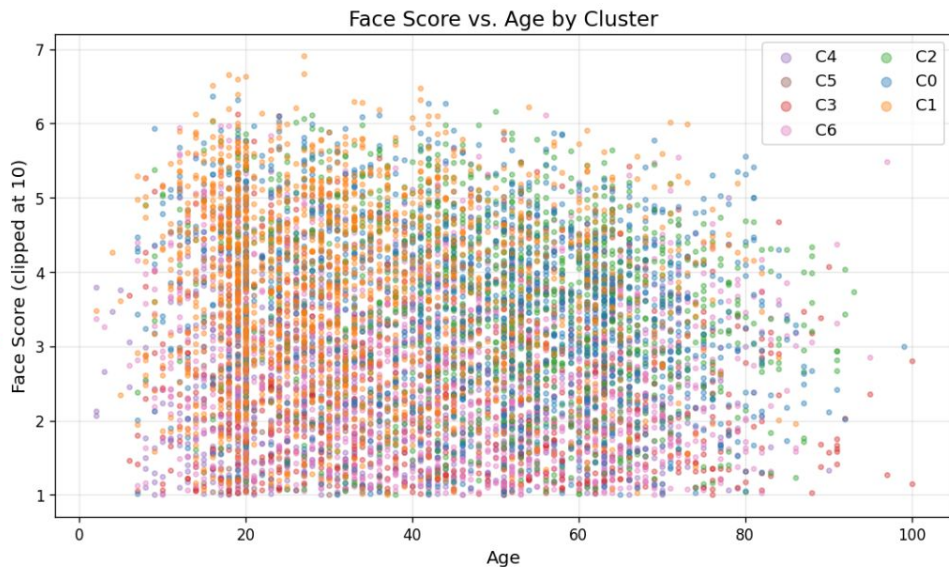
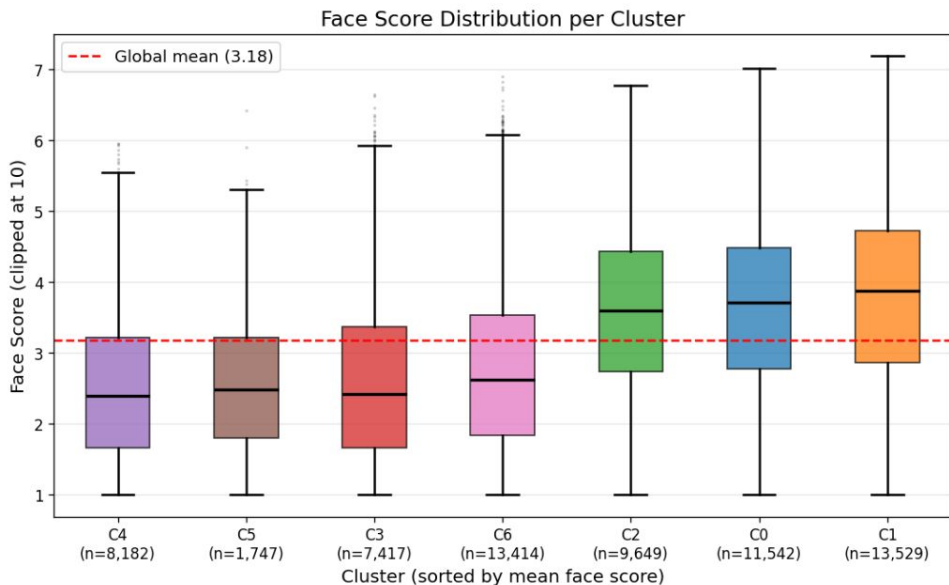
Explanation to “Images of the Same Person?”

- Cluster Purity: Fraction of the person’s images in the same cluster
 - 1: all images land in the same cluster \Rightarrow rather identity-based clustering
 - $\rightarrow 0$: all images in different clusters \Rightarrow rather age-based clustering

- Clustering vs. Age-Span:
 - Flat line: clustering is not affected by how the person ages \rightarrow rather gender, image style, identity, ...
 - Negative slope: the clustering is picking up on the age-related appearance change

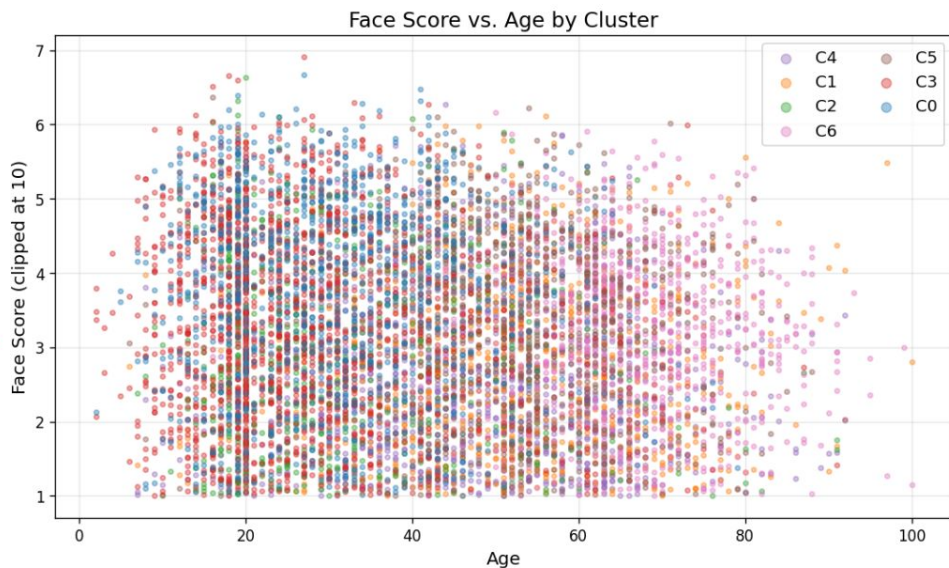
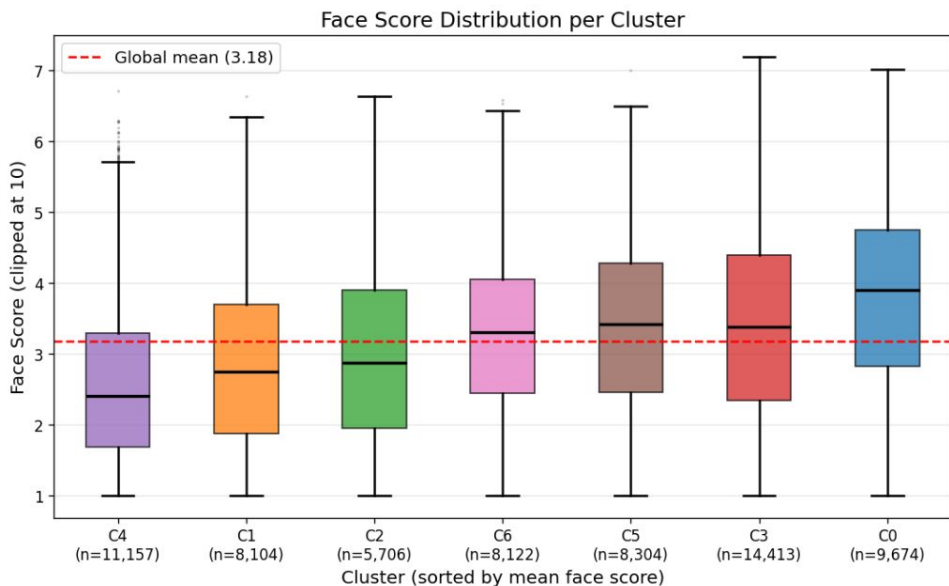
VGG16 - Facescore Bias?

Face Score Analysis - VGG16 KMeans



SAE - Facescore Bias?

Face Score Analysis - Age-Supervised AE



Explanation to “Facescore Bias?”

Left - FaceScore vs. Cluster:

- FaceScore of a cluster completely below mean: maybe grouping images that were hard to detect (blurry, side-profile, occluded) rather than by age, identity, ...
- Low number and low FaceScore might indicate a cluster-category of “bad images”

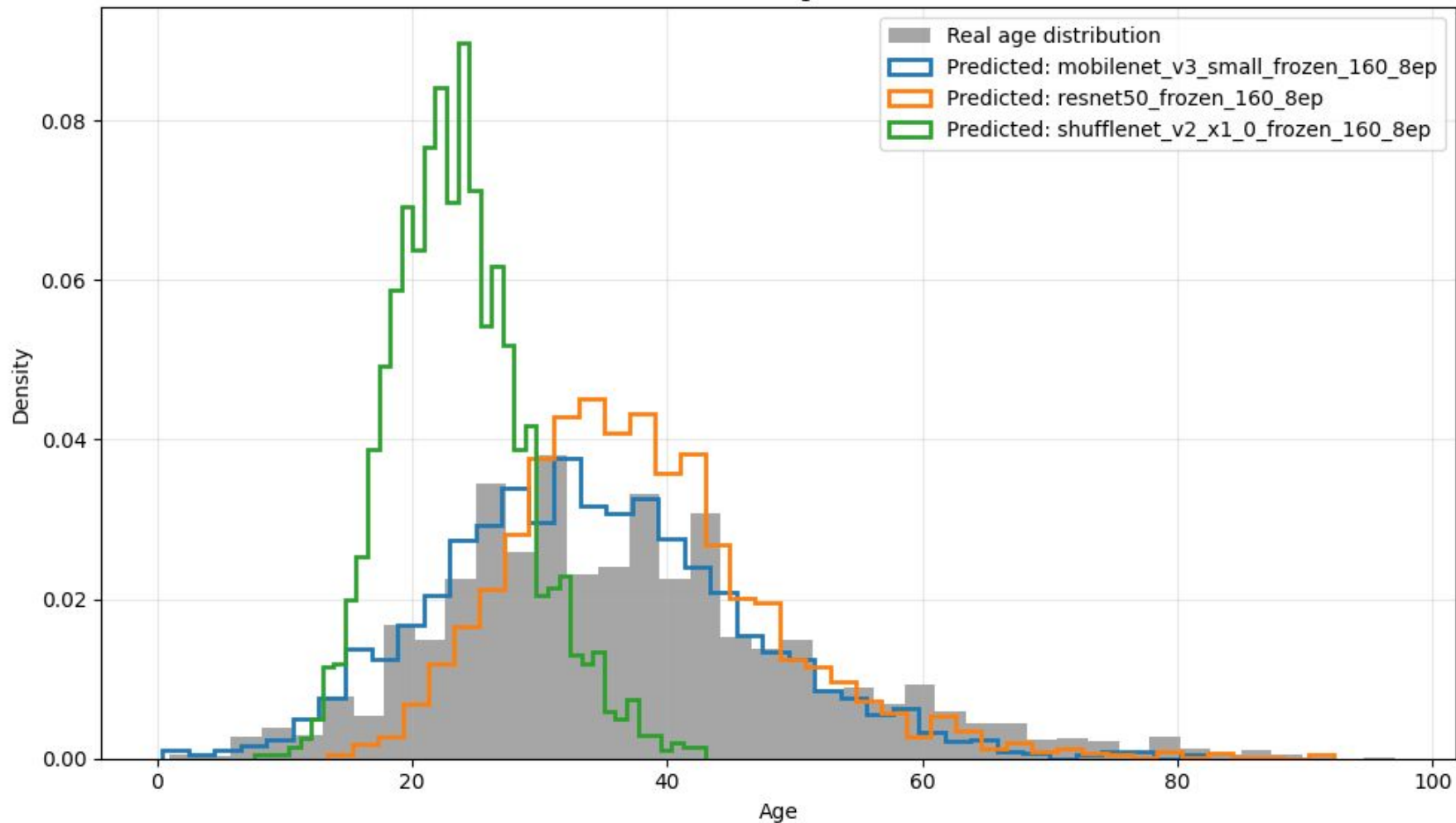
Right: Face-Score vs. Age

- Diagonal pattern (lower left triangle): might indicate that age and image quality are confounded in the cluster
- Uniform distribution: clustering rather not driven by image artifacts

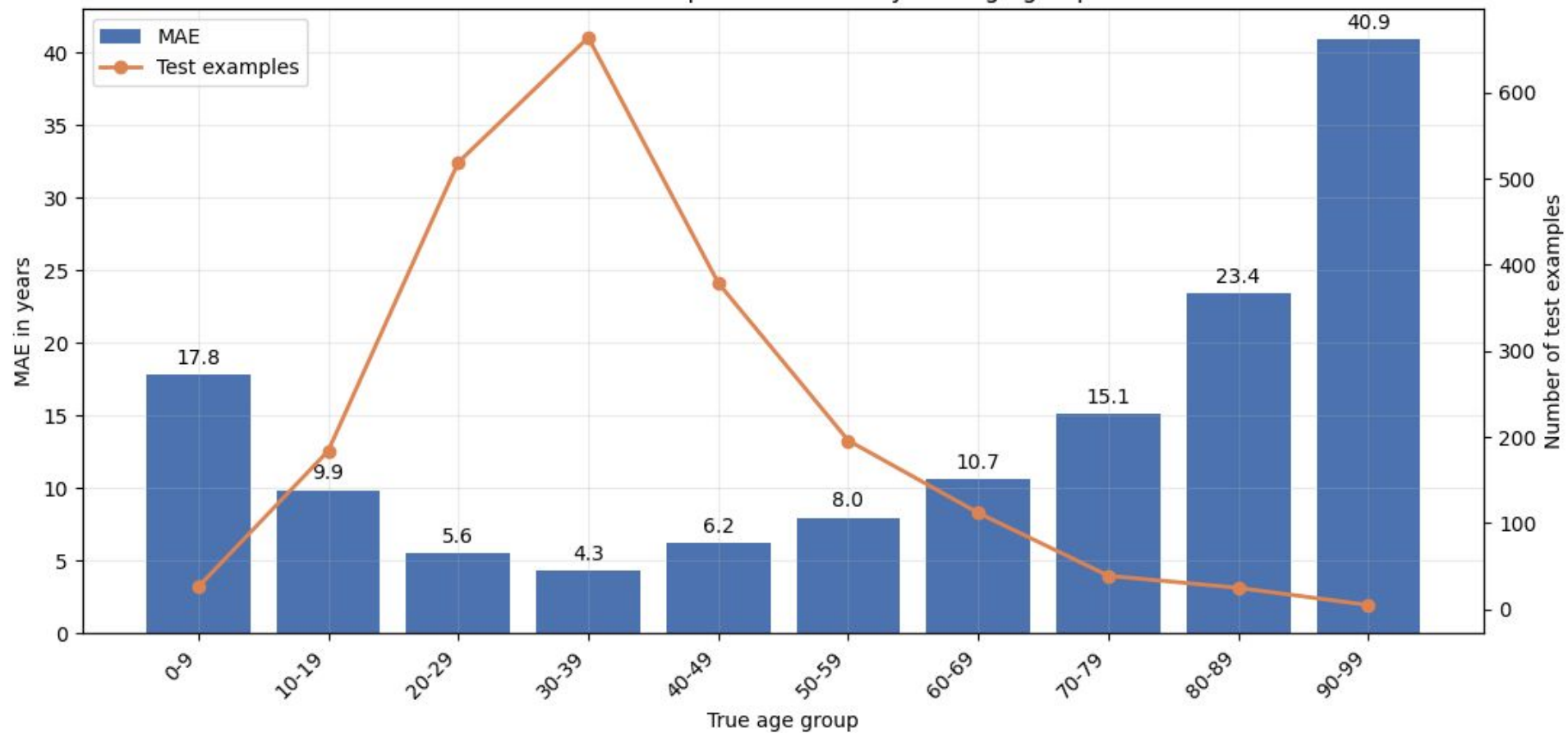
Comparison VGG16 and SAE: SAE seems to be a bit more age-driven, while the clusters C4, C5, C3 in VGG16 seem to be biased by the facescore and/or image-quality

Appendix - Regression

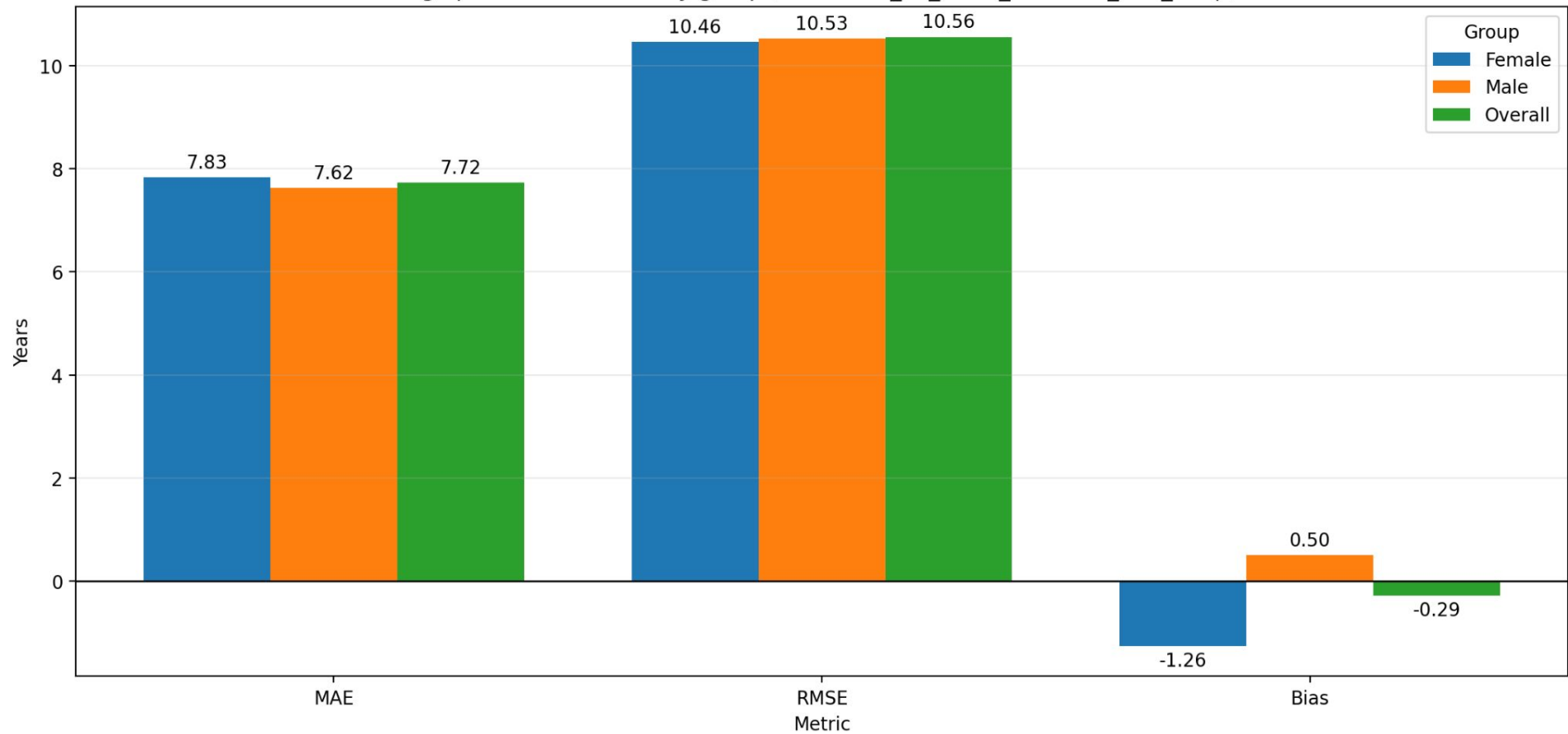
Real vs Predicted Age Distributions



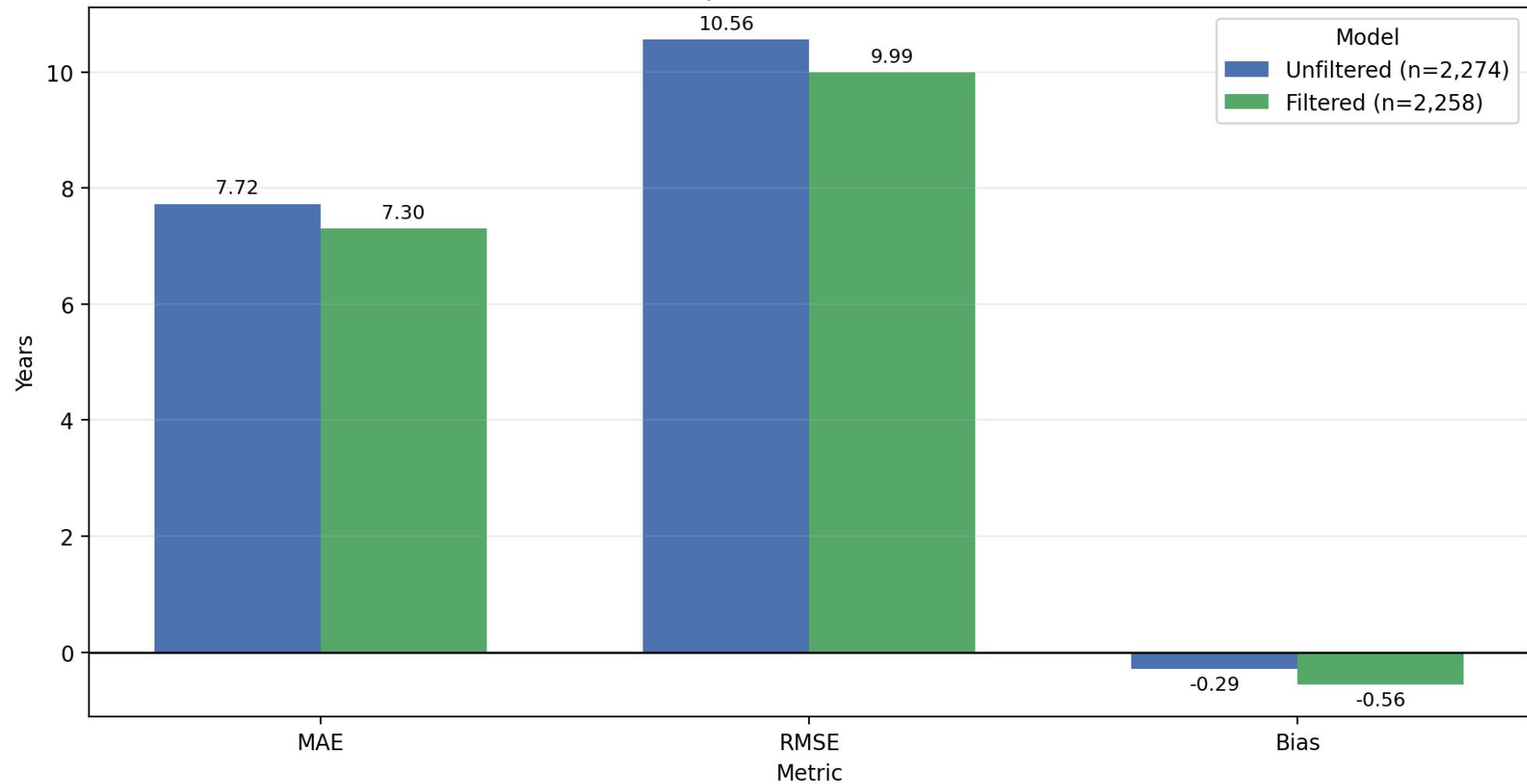
Combined model prediction MAE by true age group



Age prediction metrics by group (mobilenet_v3_small_unfrozen_224_30ep)



Overall MobileNet performance: unfiltered vs filtered



Images with Bad Predictions

true=80.0, pred=39.7, err=-40.3



true=88.0, pred=43.7, err=-44.3



true=92.0, pred=47.7, err=-44.3



true=79.0, pred=20.4, err=-58.6



true=73.0, pred=31.1, err=-41.9



true=78.0, pred=30.8, err=-47.2



true=85.0, pred=36.0, err=-49.0



true=78.0, pred=30.6, err=-47.4



Pretrained Models and ResNet50

- Training a CNN from scratch requires millions of labeled images and significant computational resources.
- Pretrained CNNs have already learned rich visual representations from large-scale datasets
- Transfer learning allows these learned features to be adapted to new tasks with much less data
- Widely used pretrained models include VGG, Inception, EfficientNet, DenseNet, and ResNet
- REsNEt50 was introduced 2015 by Microsoft

ResNet50 Architecture

- Initial 7 x 7 convolutional layer followed by max pooling
- 4 sequential residual stages with increasing feature depth
- Built using bottleneck residual blocks
 - 1 x 1 convolution (dimensionality reduction)
 - 3 x 3 convolution (feature extraction)
 - 1 x 1 convolution (feature expansion)
- Skip connections add the input directly to the block output
- Ends with Global Average Pooling followed by a Fully Connected classification layer
- Total depth: 50 layers and 25 million trainable parameters

ResNet50 Training

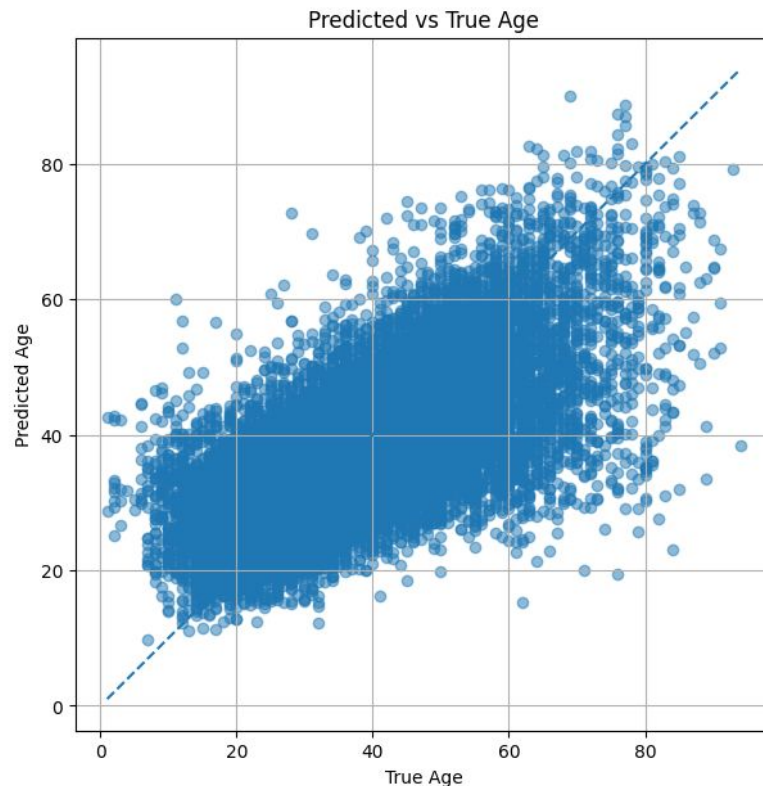
- Originally trained on the ImageNet dataset
- Imagenet contains over 1 million labeled images
- Covers 1.000 objects and categories
- Training uses supervised learning with stochastic gradient descent (SGD)
- Data augmentation techniques improve robustness and generalization

How Did We Apply ResNet50

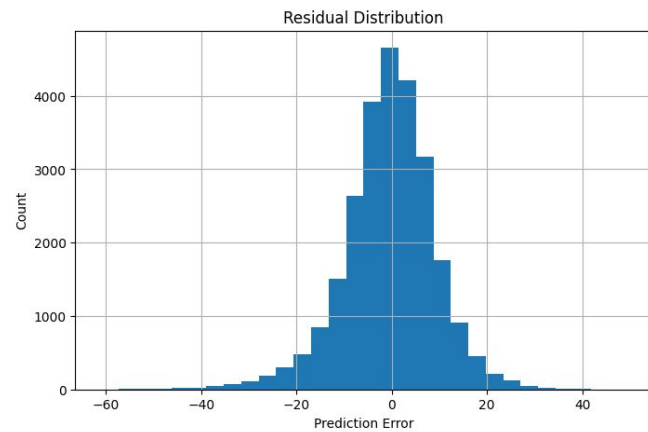
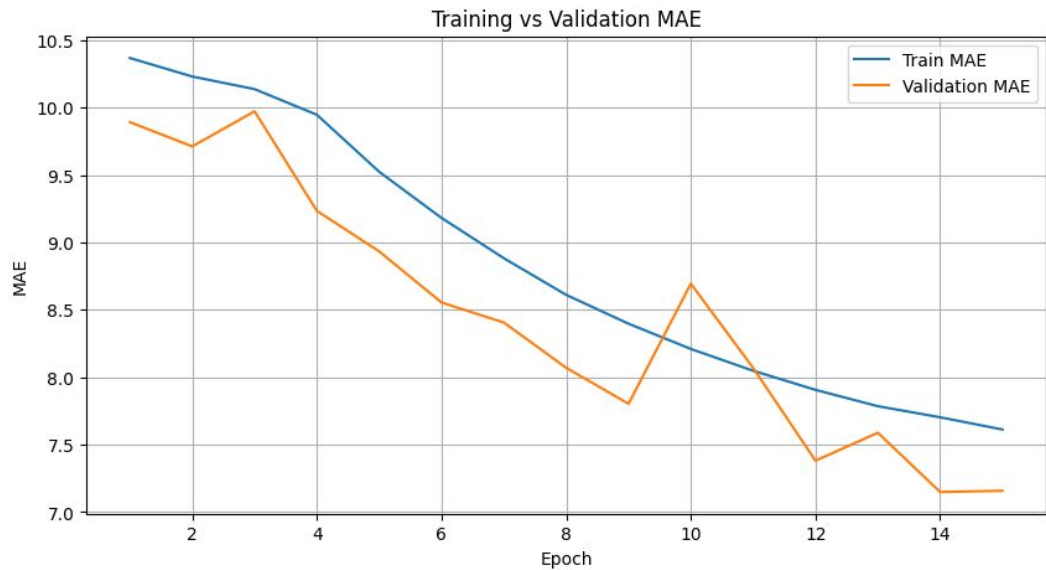
- Initialize the network using pretrained weights instead of random initialization
- Replace the final classification layer with task-specific output layers (Age Regression and Gender Classification)
- Freeze early convolutional layers to preserve general visual features
- Fine-tune deeper layers to learn facial characteristics relevant to age and gender

CNN from Scratch Architecture

- 4 convolutional blocks
- 8 convolutional layers
- Batch Normalization after every convolution
- ReLU activation through out
- Feature expansion: $32 \rightarrow 64 \rightarrow 128 \rightarrow 256$
- Parameters

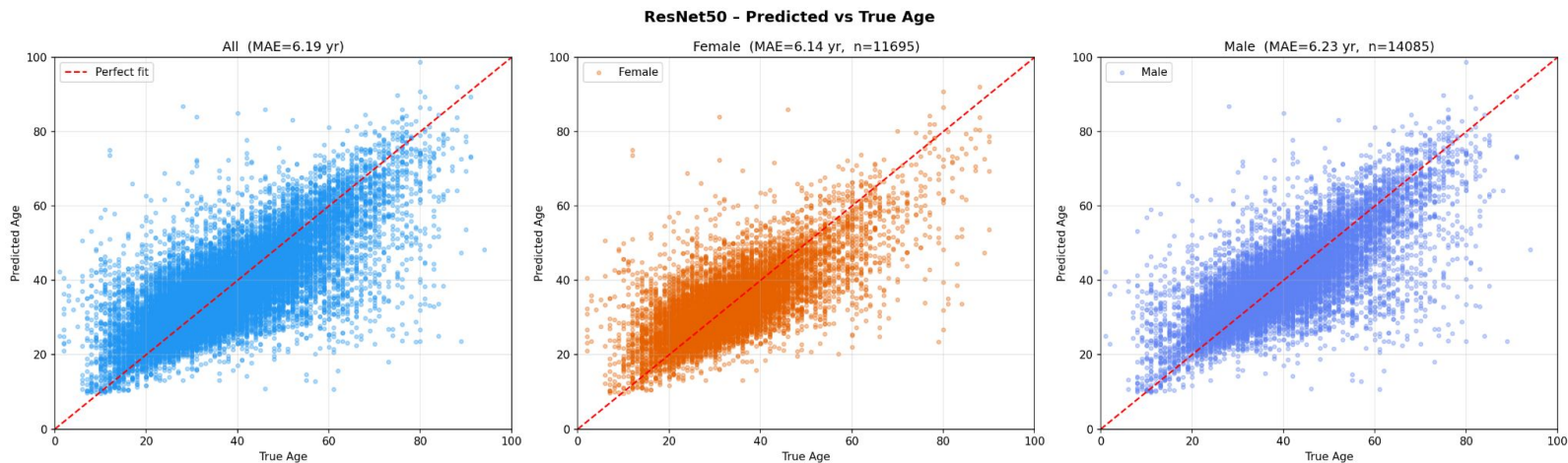


CNN for Age Regression from Scratch



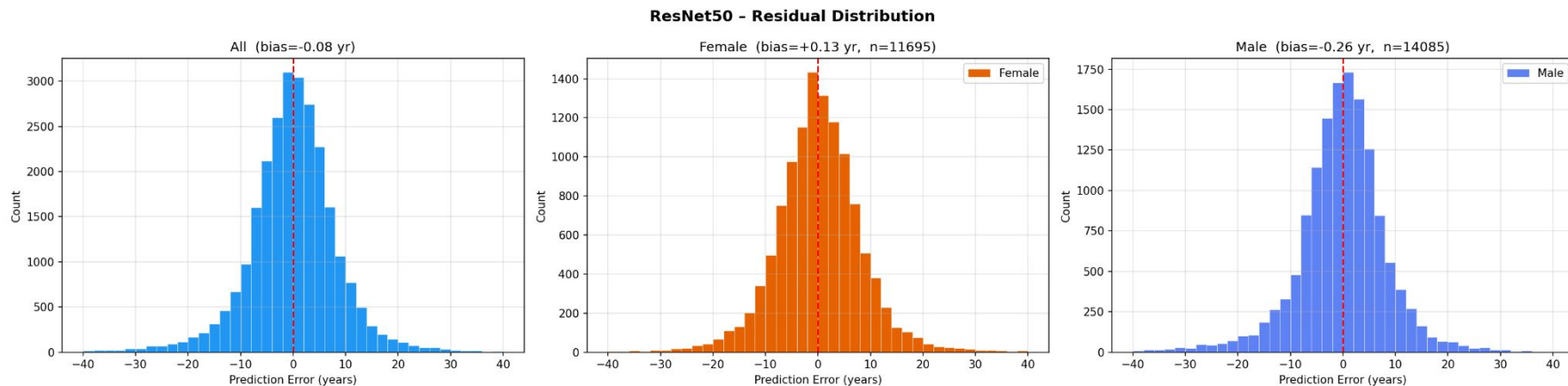
Age Regression ResNet50

True age vs predicted age by gender



Basic ResNet50 Regression Split by Gender

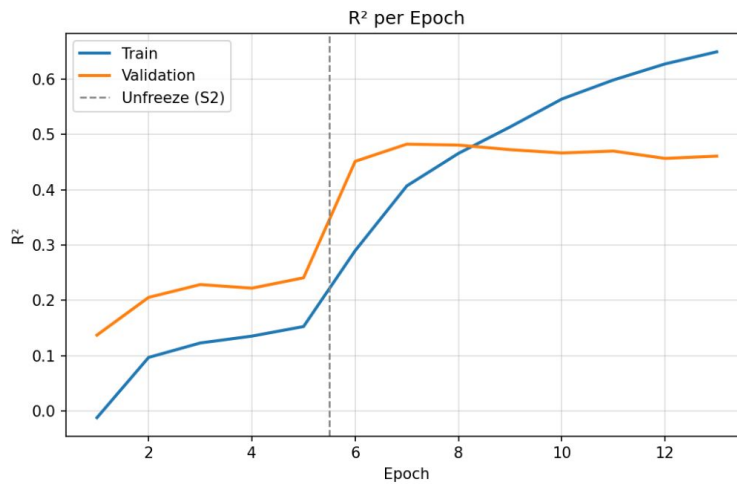
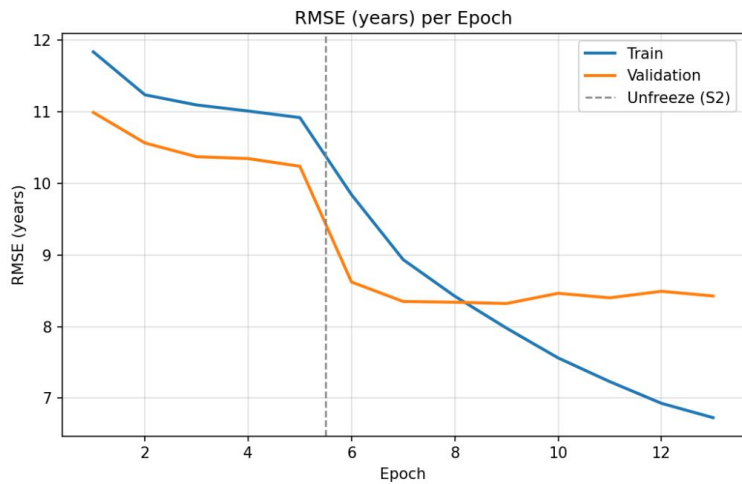
Residuals by Gender



Age Regression ResNet50

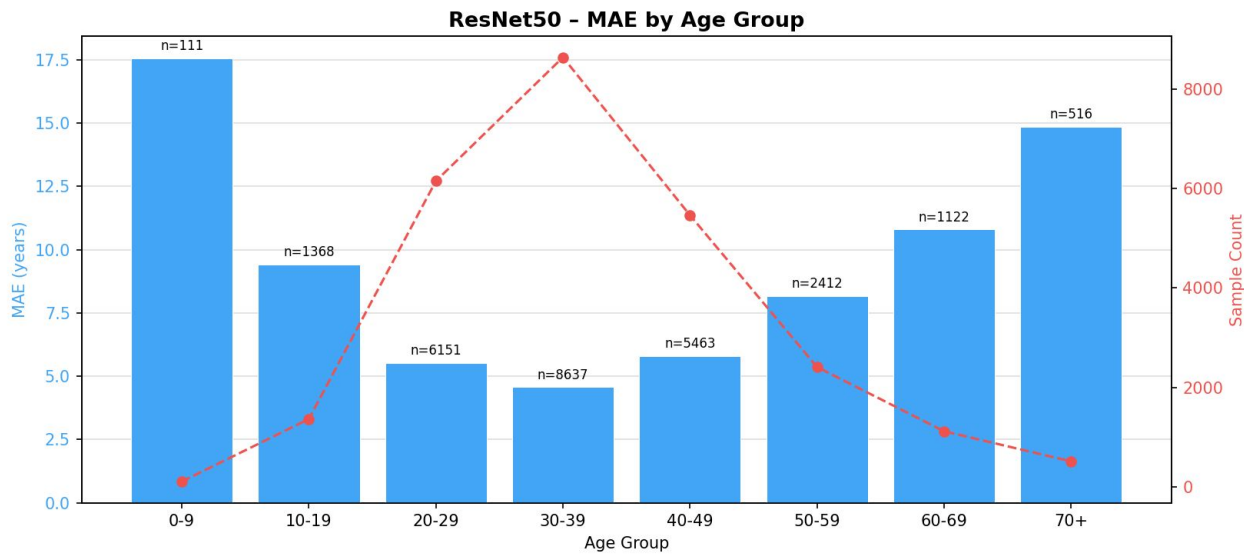
Some more metrics, show the same behaviour as MAE

ResNet50 - RMSE & R² Curves



Age Regression Using ResNet50

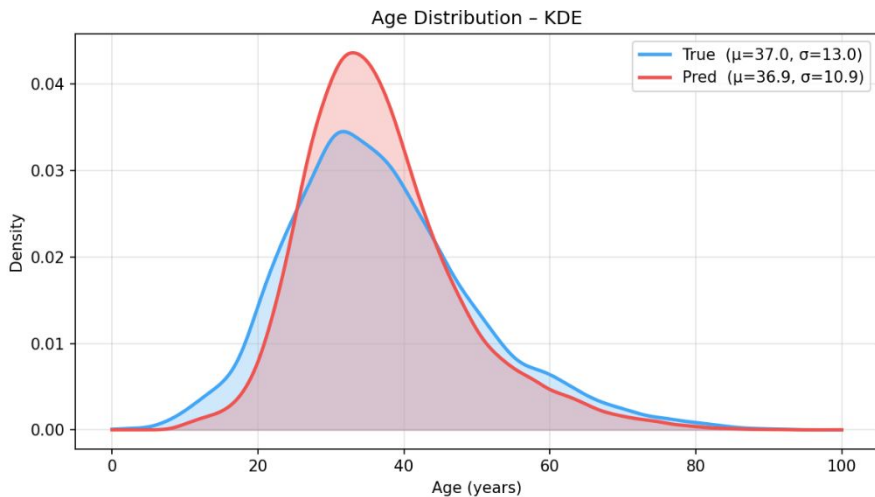
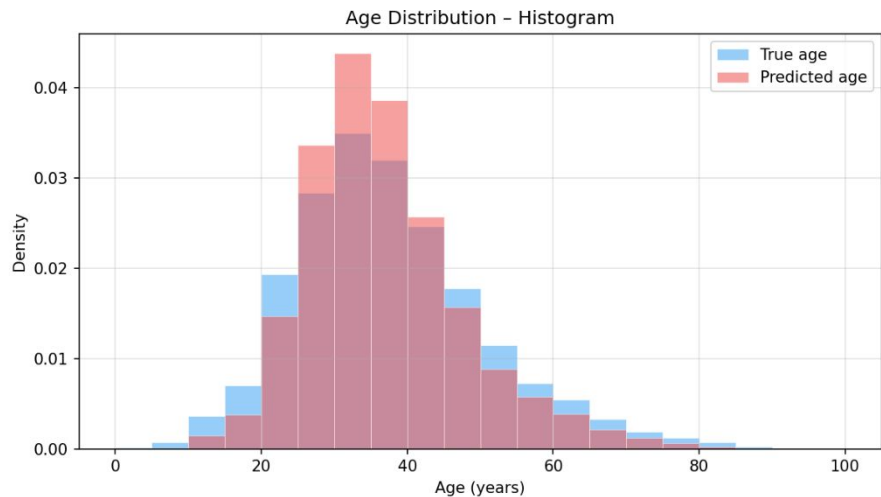
MAE compared to the number of samples in this age group (inverse correlation)



Age Regression Using ResNet50

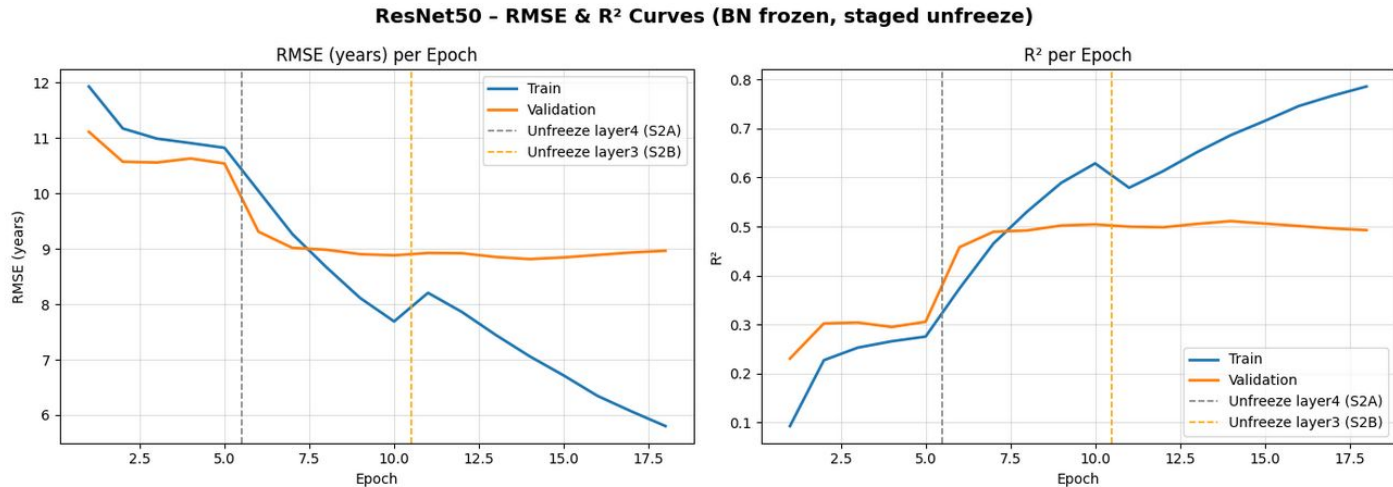
Predicted vs actual distribution of ages. Left out for time reasons

ResNet50 - True vs Predicted Age Distribution



Age Regression Using ResNet50 Improved

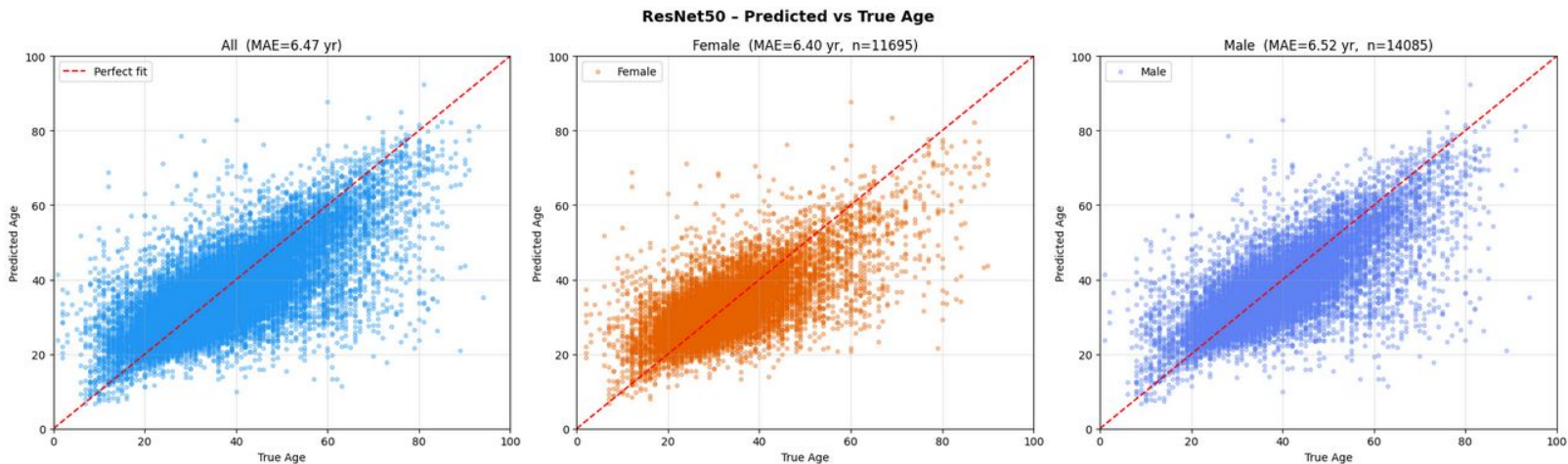
Some more metrics, show the same behaviour as MAE (multiple stage case)



Saved: ResNet50_rmse_r2_curves.png

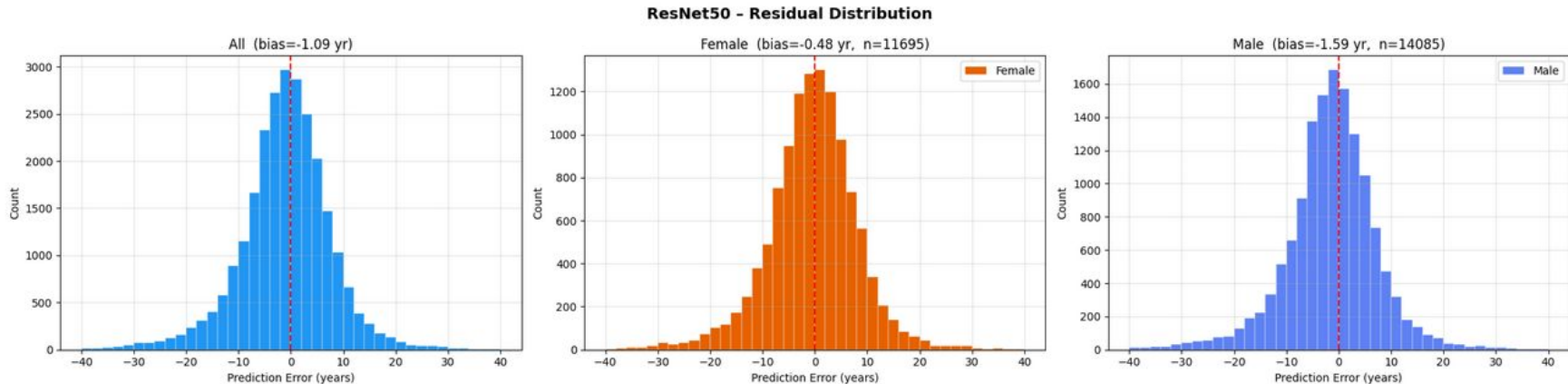
Age Regression ResNet 50 Improved

True age vs predicted age for “improved” unfreezing by gender



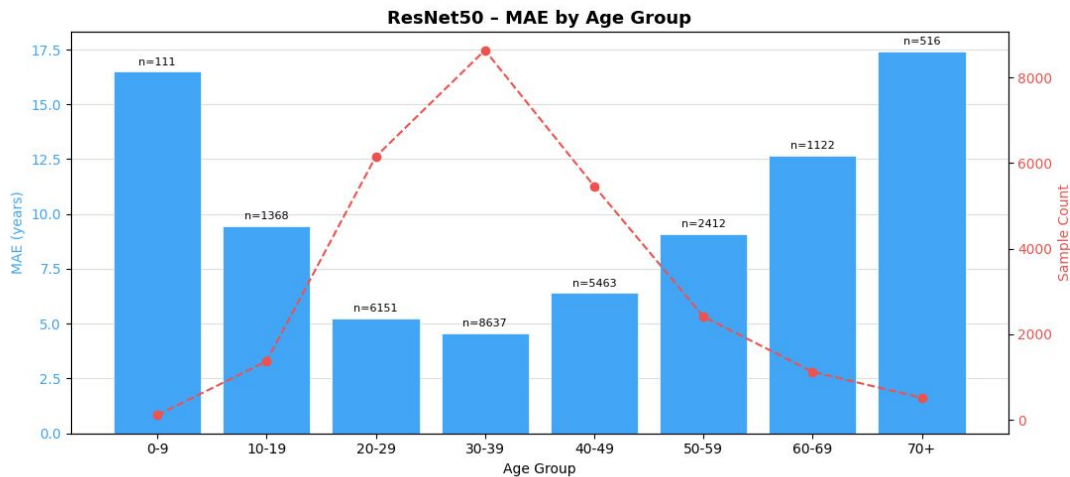
Age Regression ResNet 50 Improved

Residuals for “improved” unfreezing by gender



Age Regression ResNet 50 Improved

MAE compared to the number of samples in this age group (inverse correlation) for “improved” unfreezing

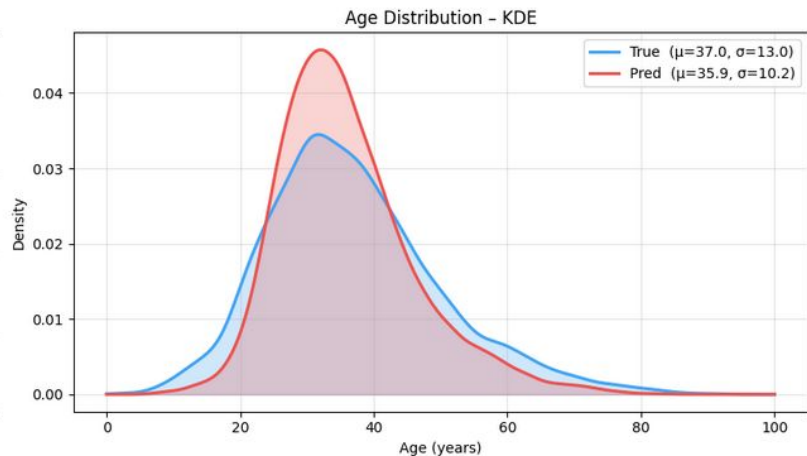
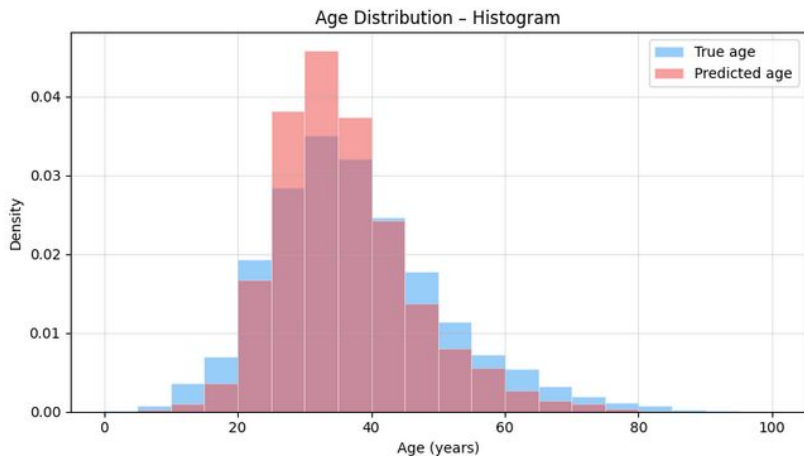


Saved: ResNet50_mae_by_age_group.png

Age Regression ResNet 50 Improved

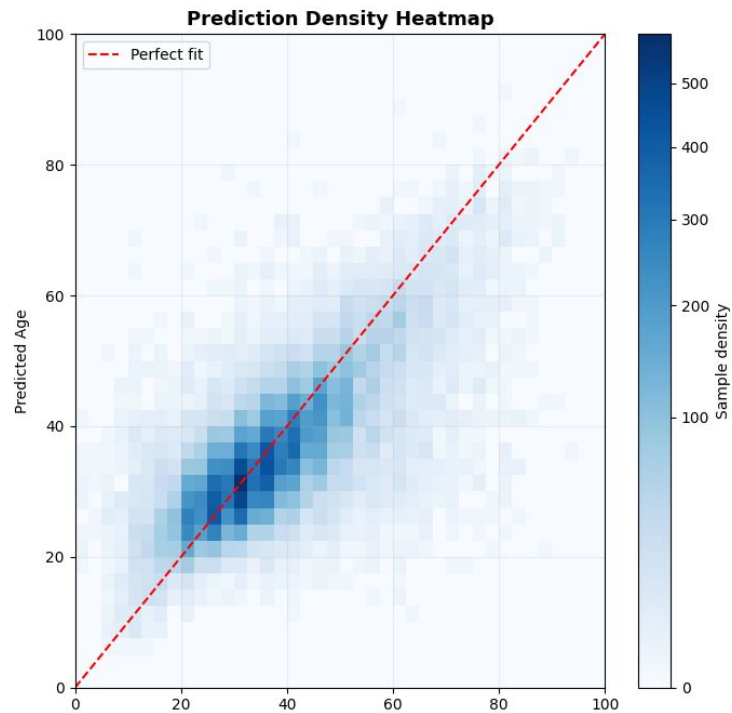
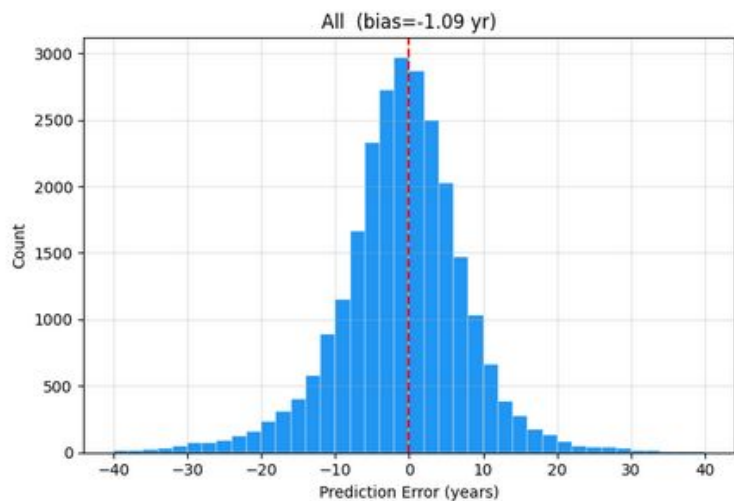
Predicted vs actual distribution of ages for “improved” unfreezing. Left out for time reasons

ResNet50 - True vs Predicted Age Distribution



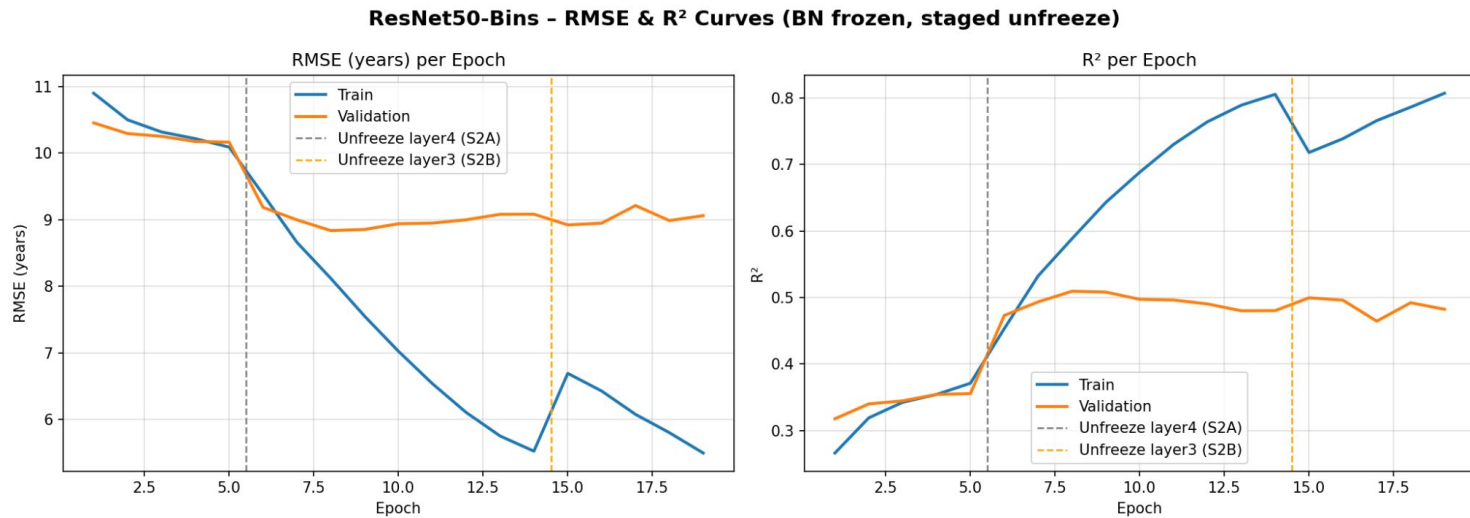
Appendix - Age Classification

Age Binning



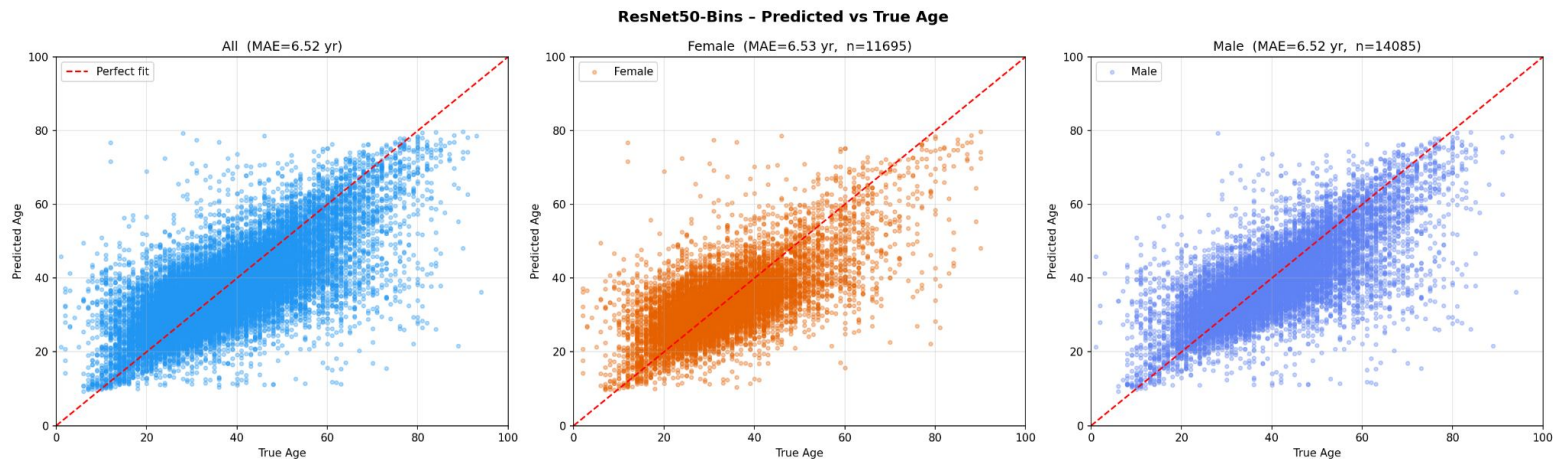
Age Binning

Some additional metrics



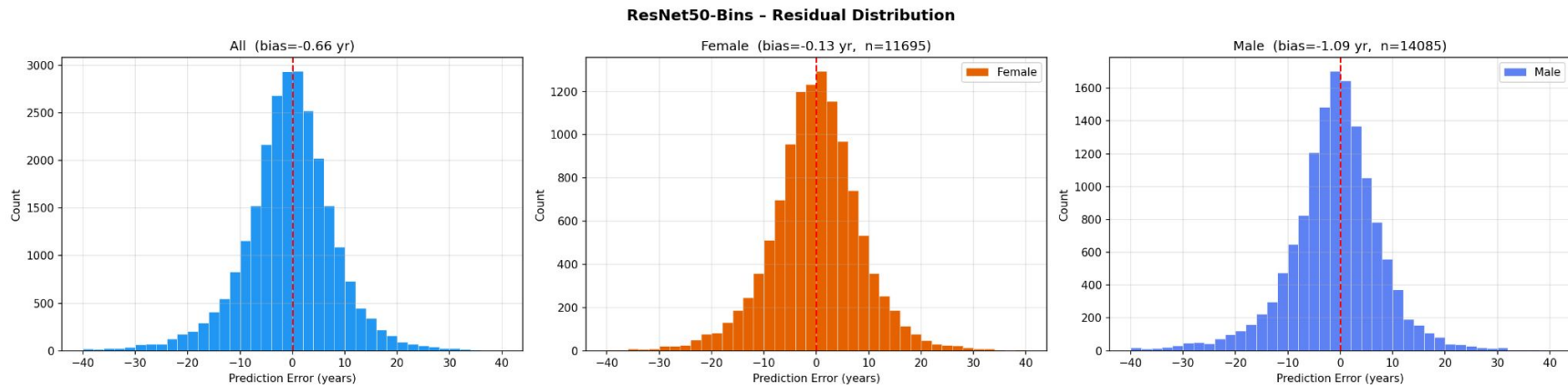
Age Binning

True vs predicted age



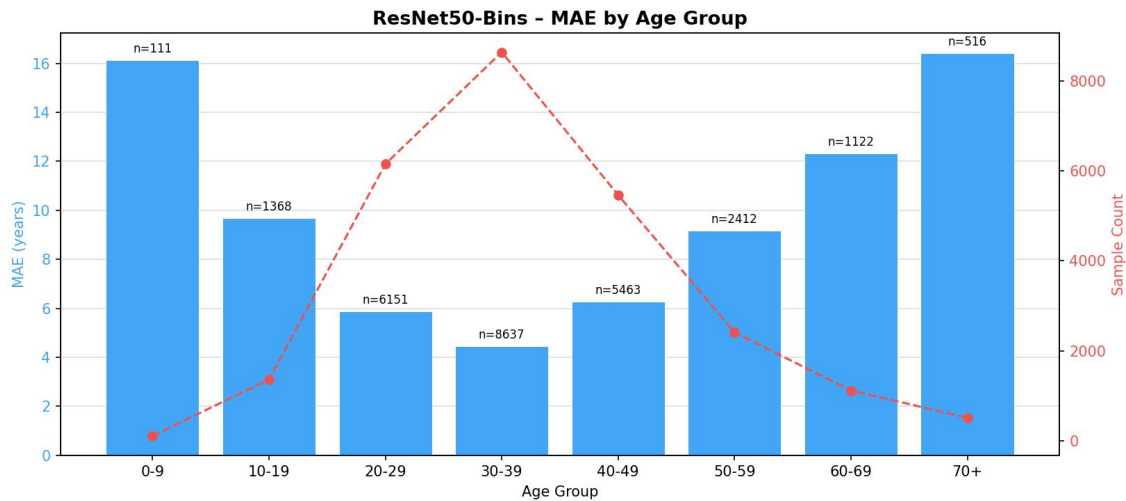
Age Binning

Residuals by gender for the age binning setting



Age Binning

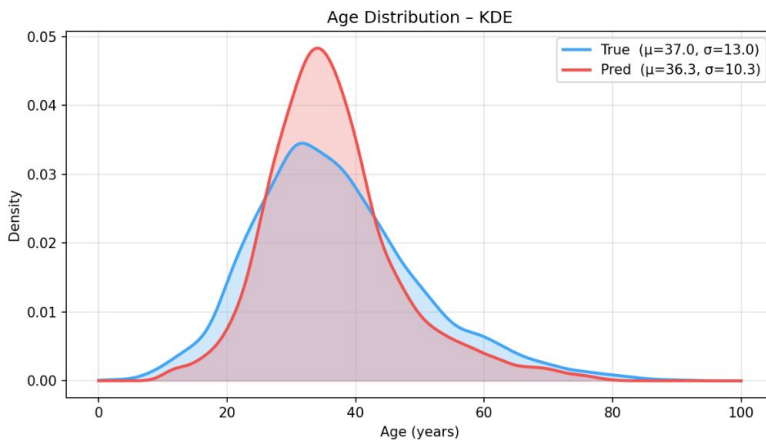
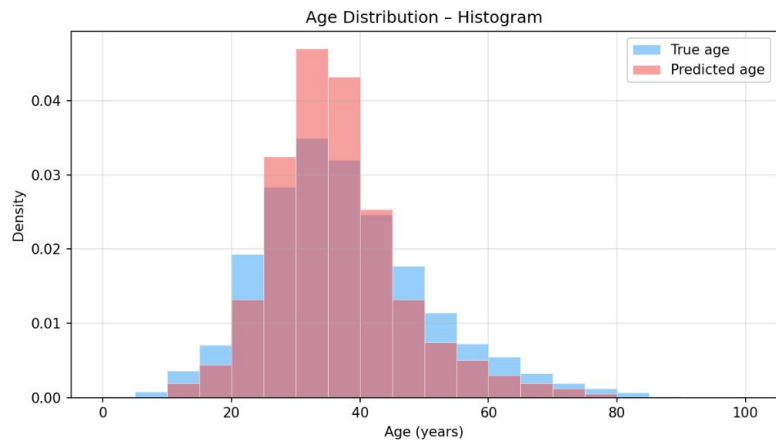
MAE compared to the number of sample for the classification approach



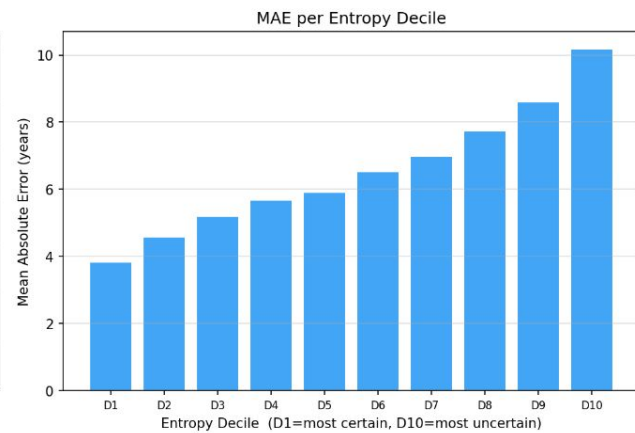
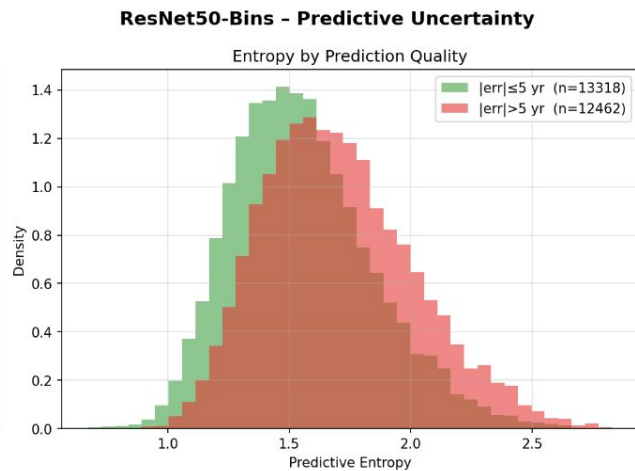
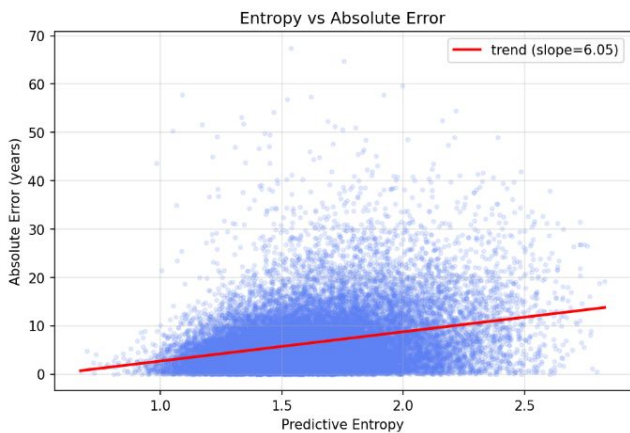
Age Binning

True vs predicted age distribution by the classification approach

ResNet50-Bins - True vs Predicted Age Distribution

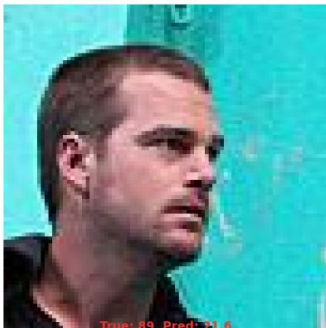


Age Binning



Best and Worst Age Predictions

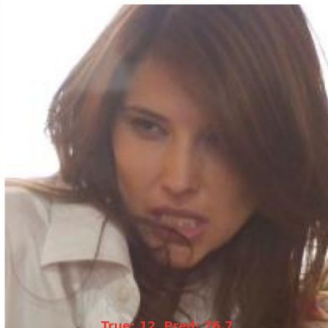
True: 39 Pred: 39.0
|err|=0.0 yr



True: 89 Pred: 21.6
|err|=67.4 yr



True: 29 Pred: 29.0
|err|=0.0 yr



True: 12 Pred: 76.7
|err|=64.7 yr



True: 34 Pred: 34.0
|err|=0.0 yr



True: 12 Pred: 71.4
|err|=59.6 yr



True: 32 Pred: 32.0
|err|=0.0 yr



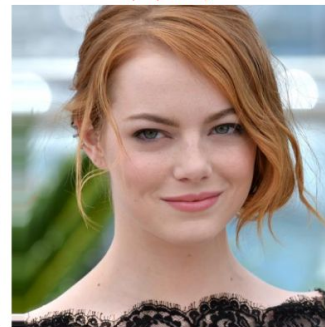
True: 59 Pred: 26.2
|err|=57.8 yr



True: 24 Pred: 24.0
|err|=0.0 yr



True: 84 Pred: 26.2
|err|=57.8 yr



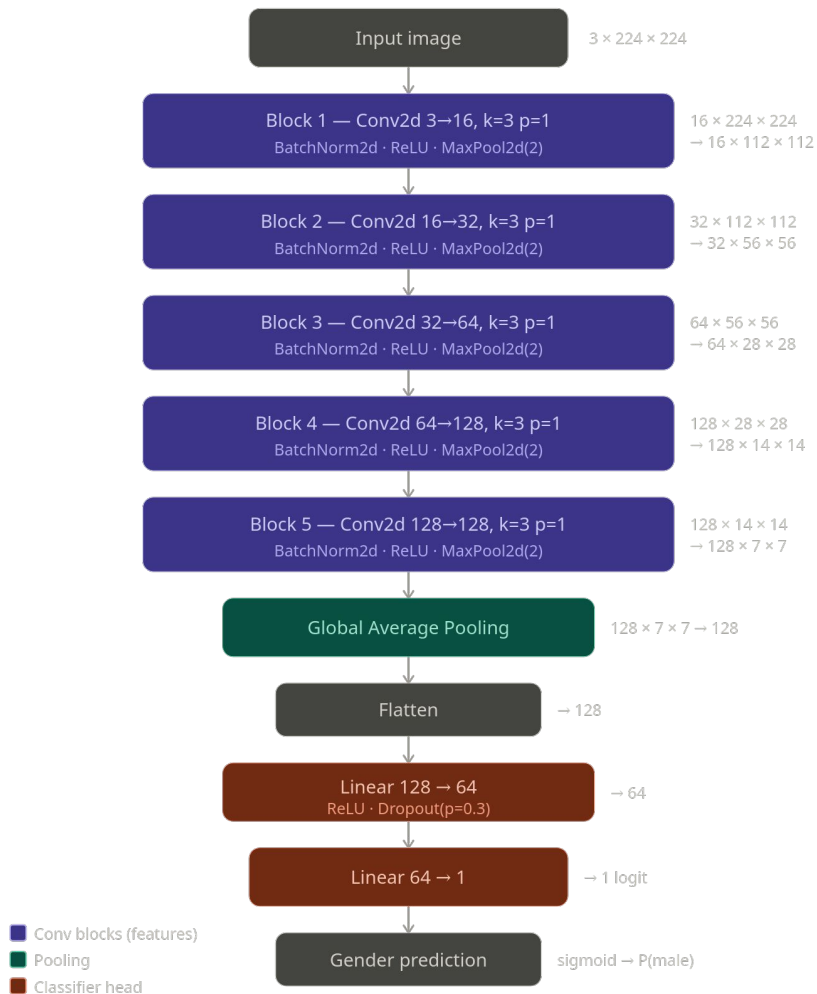
Appendix - Gender Classification

Baseline CNN

Runtime

▶ 4h 57m 31s · GPU T4 ×2

total params: 253,713



ResNet50 and MobileNetV3

ResNet50 and MobileNetV3 for full dataset:

Runtime

▶ 3h 59m 25s · GPU T4 ×2

```
ResNet-50      total: 23,510,081, trainable: 23,510,081
MobileNetV3    total: 4,203,313 trainable: 4,203,313
```

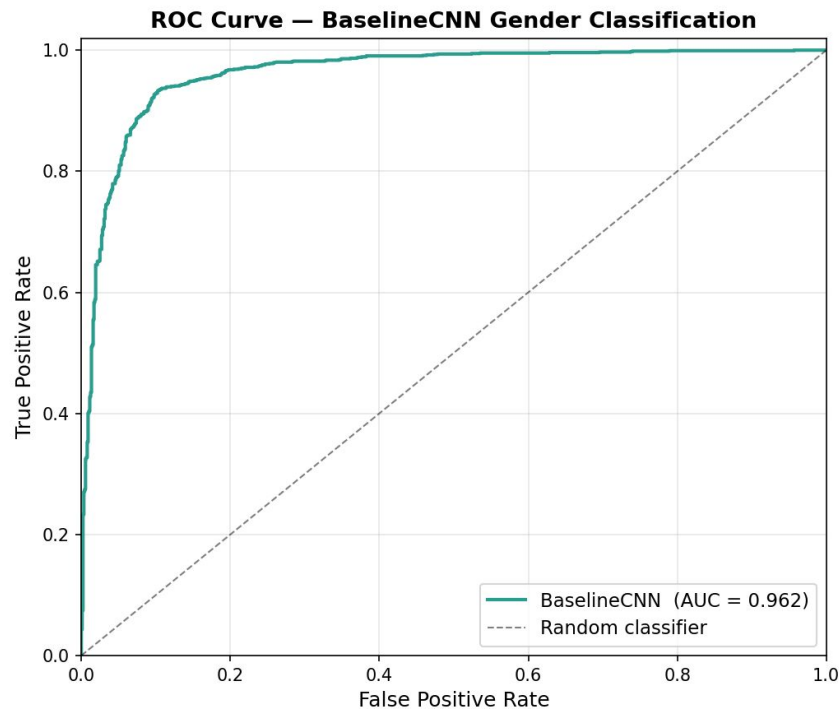
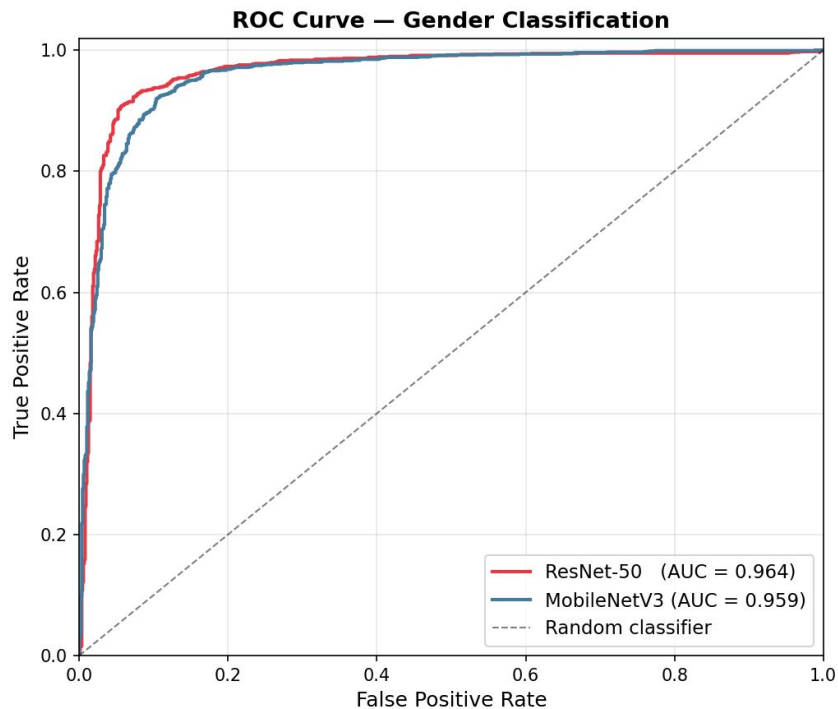
ResNet50 and MobileNetV3 for 1 pic per celeb dataset:

Runtime

▶ 31m 5s · GPU T4 ×2

```
ResNet-50      total: 23,510,081 trainable: 23,510,081
MobileNetV3    total: 4,203,313 trainable: 4,203,313
```

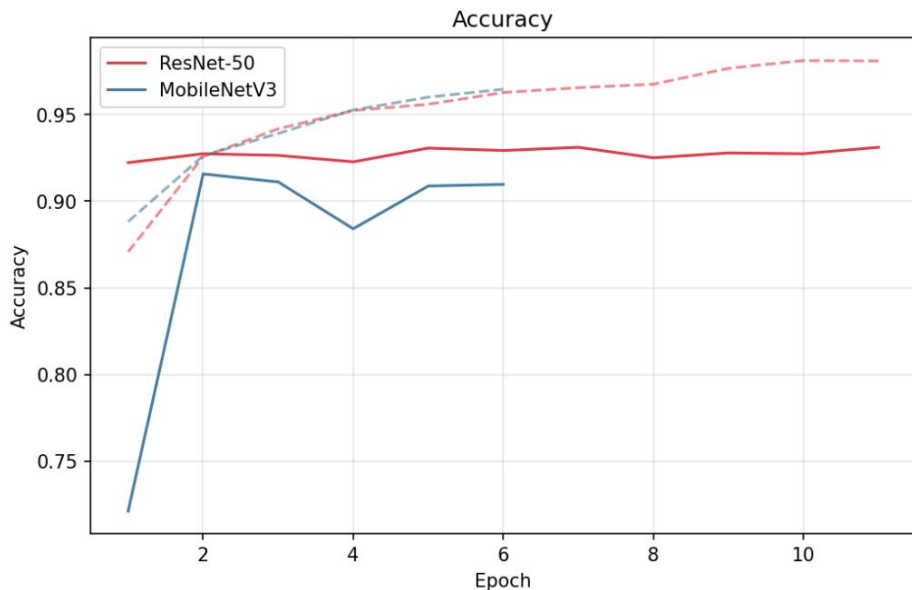
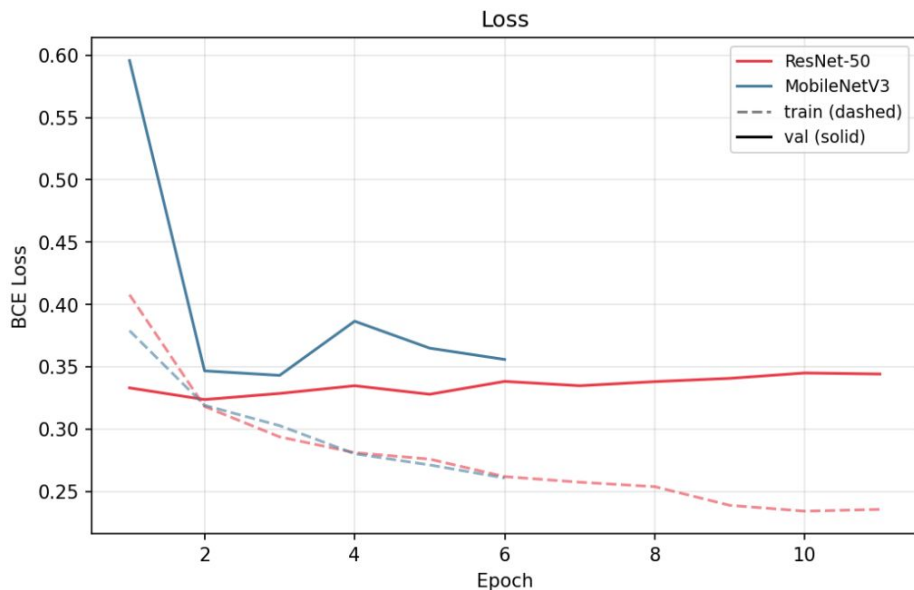
Gender Classification - Model Evaluation



dataset: best face score picture by celeb ~15000 pictures

Gender Classification - Model Evaluation

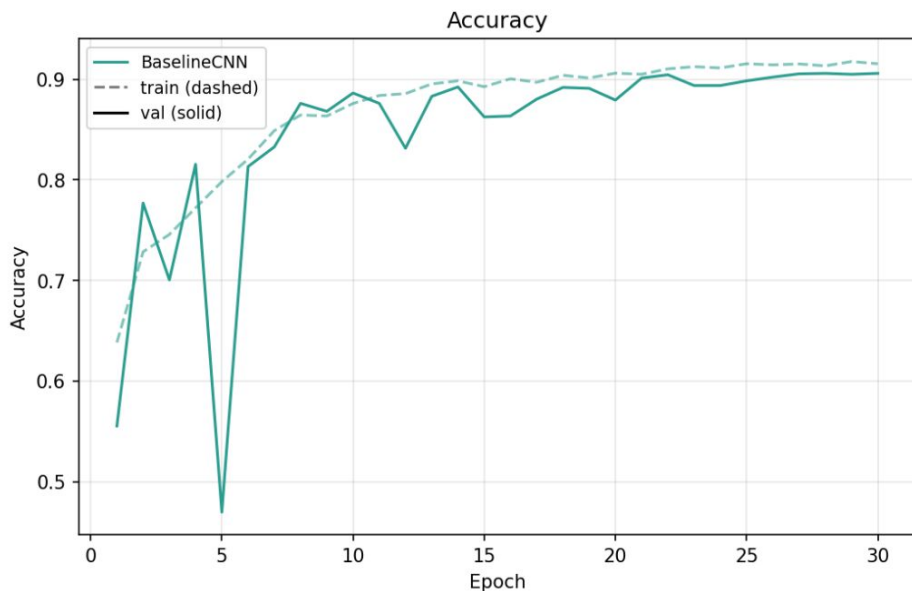
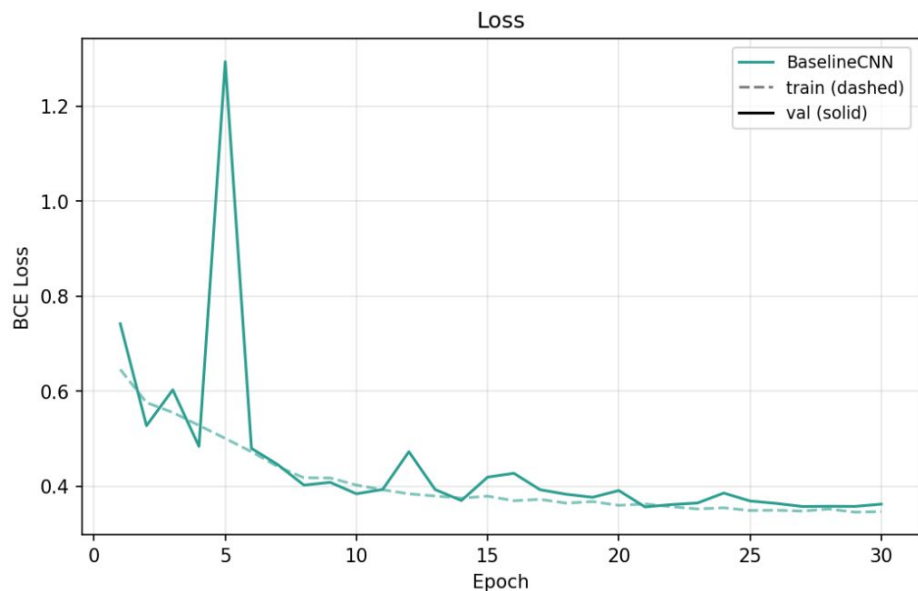
Training History



dataset: best face score picture by celeb ~15000 pictures

Gender Classification - Model Evaluation

BaselineCNN — Training History

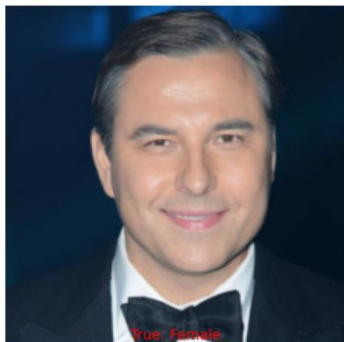


dataset: best face score picture by celeb ~15000 pictures

Gender Classification - mobilenet evaluation

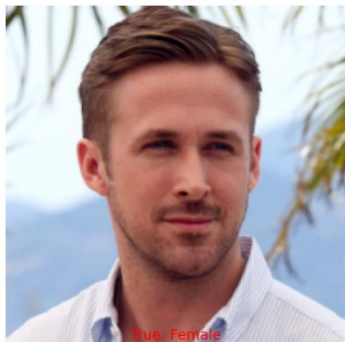
mobilenetv3 — Top-5 Confident Correct vs Top-5 Confident Wrong

True: Male
Pred: Male (97%)



True: Female
Pred: Male (96%)

True: Male
Pred: Male (97%)



True: Female
Pred: Male (96%)

True: Male
Pred: Male (97%)



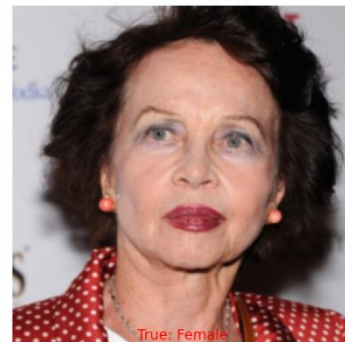
True: Male
Pred: Female (95%)

True: Male
Pred: Male (97%)



True: Female
Pred: Male (95%)

True: Female
Pred: Female (97%)



True: Female
Pred: Male (94%)



Gender Classification - Baseline CNN Evaluation

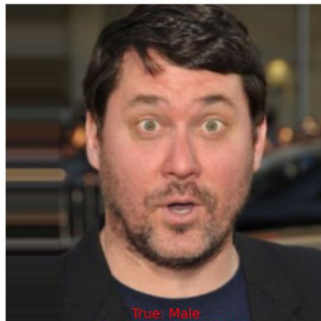
baseline_cnn — Top-5 Confident Correct vs Top-5 Confident Wrong

True: Male
Pred: Male (99%)



True: Female
Pred: Male (98%)

True: Male
Pred: Male (99%)



True: Male
Pred: Female (98%)

True: Male
Pred: Male (99%)



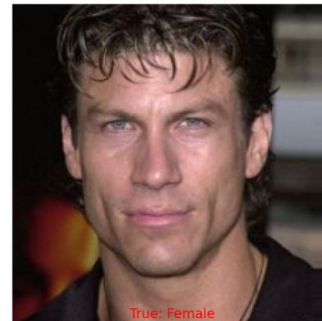
True: Female
Pred: Male (98%)

True: Male
Pred: Male (99%)

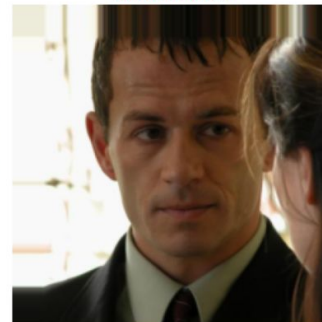
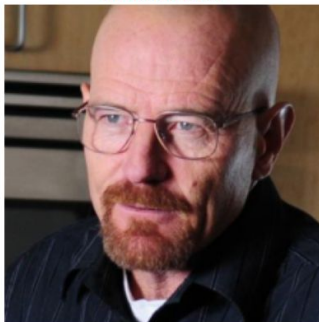
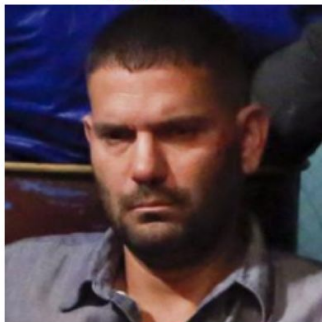


True: Female
Pred: Male (98%)

True: Male
Pred: Male (99%)



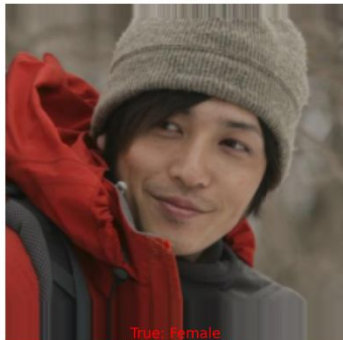
True: Female
Pred: Male (98%)



Gender Classification - ResNet Evaluation

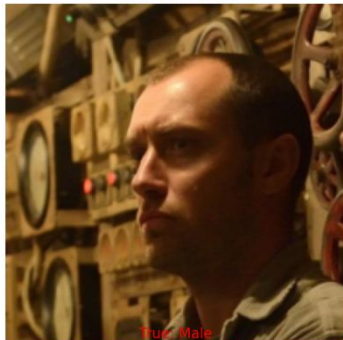
resnet50 — Top-5 Confident Correct vs Top-5 Confident Wrong

True: Male
Pred: Male (99%)



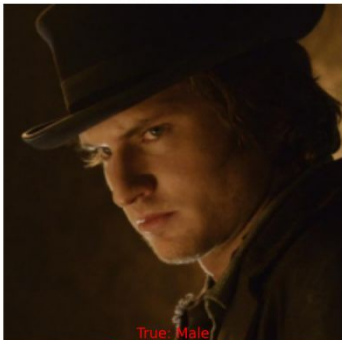
True: female
Pred: Male (98%)

True: Male
Pred: Male (99%)



True: Male
Pred: Female (98%)

True: Male
Pred: Male (99%)



True: Male
Pred: Female (98%)

True: Male
Pred: Male (98%)

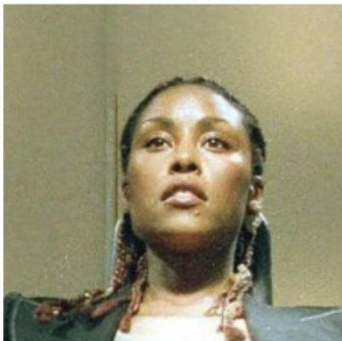


True: Male
Pred: Female (98%)

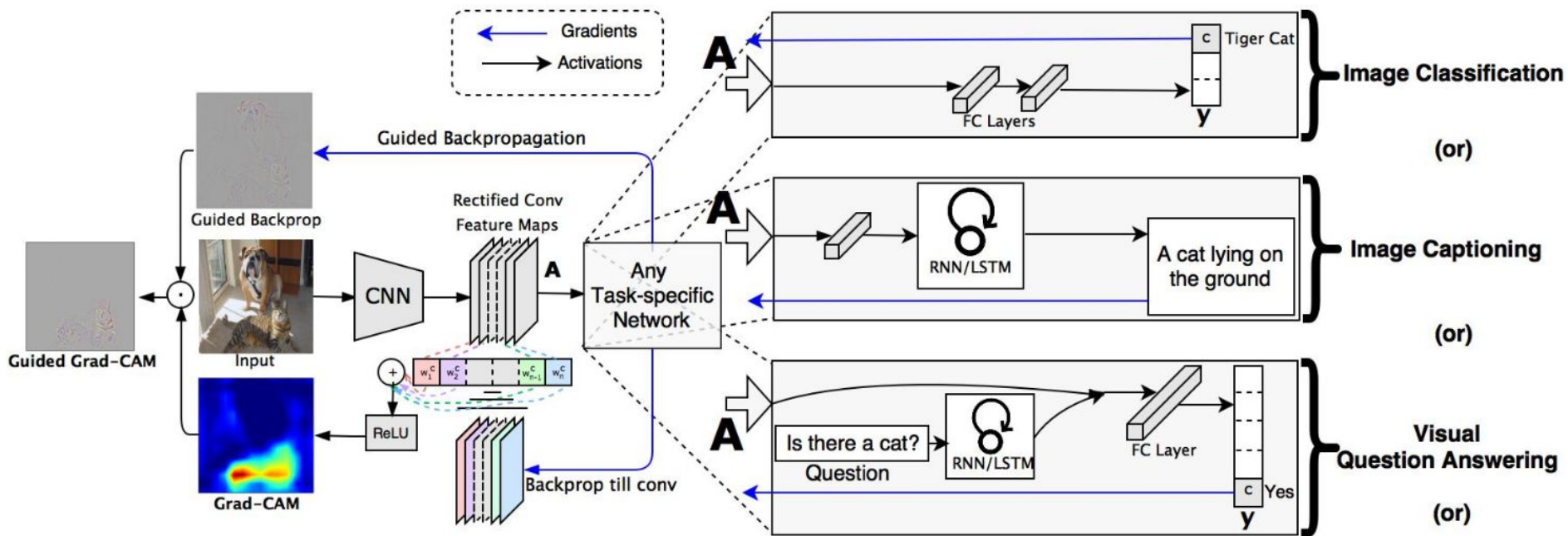
True: Male
Pred: Male (98%)

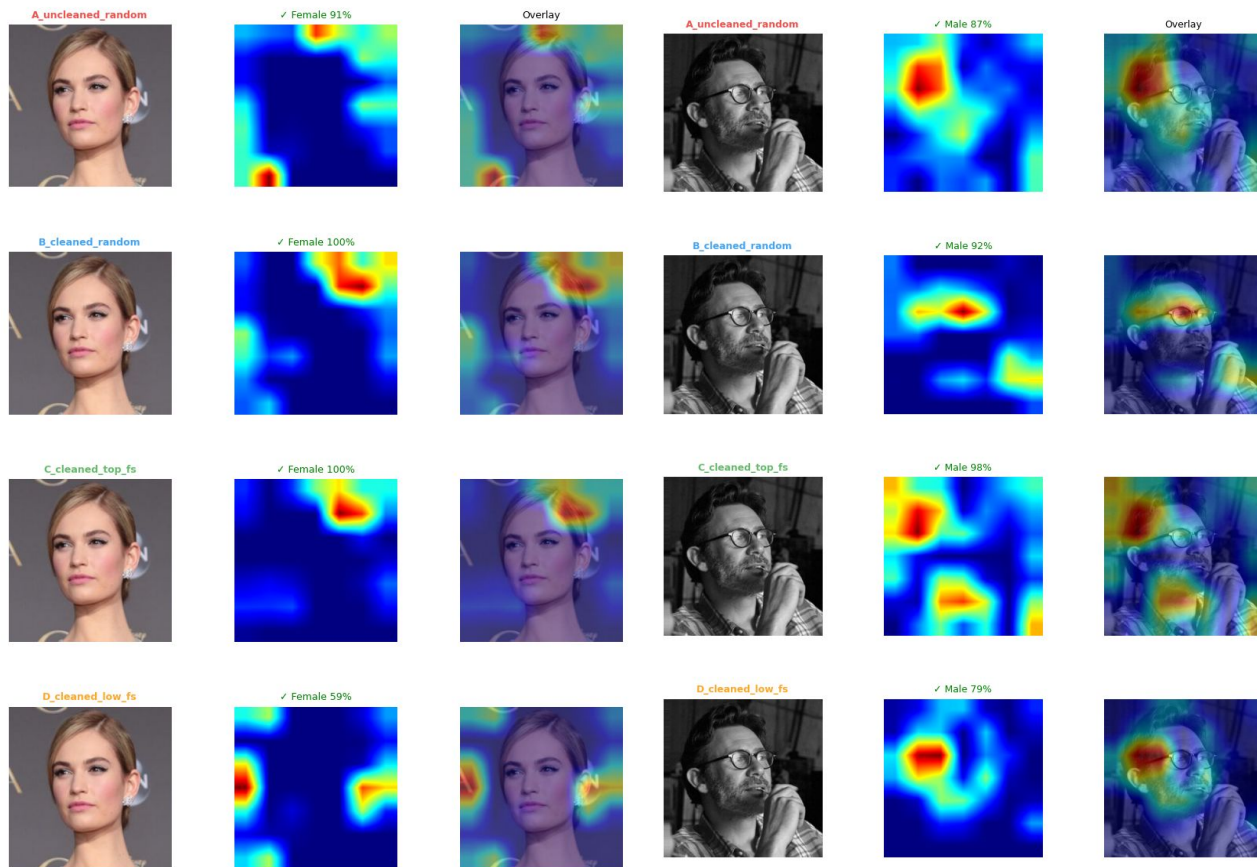


True: Female
Pred: Male (98%)



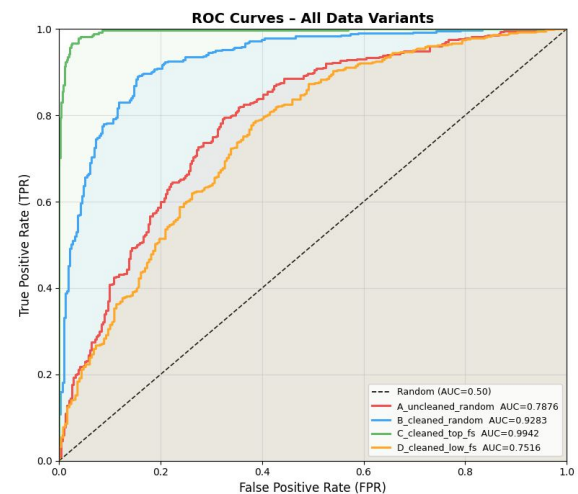
Grad-CAM Architecture





referring to slide 12 models

light mobilenet



Gender Classification - Improving Noisy Labels

Procedure:

- Sort out confidently correct classified pictures by model with label mismatch
- Evaluate model again

Noise candidates — confident wrong predictions (>90%)
Label = stored ground truth | Pred = model output — model likely correct on these

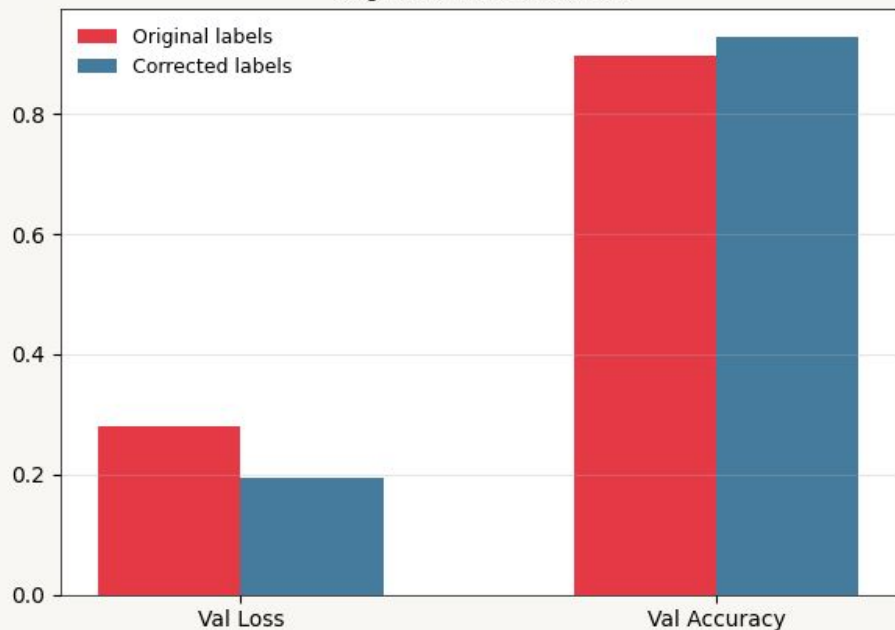


ResNet50 very confidently off predictions

Gender Classification - Improving Noisy Labels

Noise correction analysis — resnet50 (threshold=90%)

Original vs noise-corrected



ROC curve

