

THE GEORGE
WASHINGTON
UNIVERSITY
WASHINGTON, DC

AN INTERPRETABLE MACHINE LEARNING MODEL FOR GENDER PREDICTION USING SHAP AND CLASS ACTIVATION MAPS

Capstone Project

- Renzo P. Castagnino



Outline

1. Introduction

- Some News
- Trend and Challenges

2. Main question

3. Data collection & Preprocessing

4. Training models

1. VGG-16
2. ResNet
3. Inception v3

5. Interpretability models

1. SHAP
2. Class Activation Maps

6. Results

7. Conclusion

8. Future Research

Let's Begin with
some news!

One Month, 500,000 Face Scans: How China Is Using A.I. to Profile a Minority

In a major ethical leap for the tech world, Chinese start-ups have built algorithms that the government uses to track members of a largely Muslim minority group.



SenseTime is among the Chinese artificial intelligence companies developing facial recognition technology. Gilles Sabrié for The New York Times



By Paul Mozur

April 14, 2019



阅读简体中文版 [阅读繁體中文版](#)

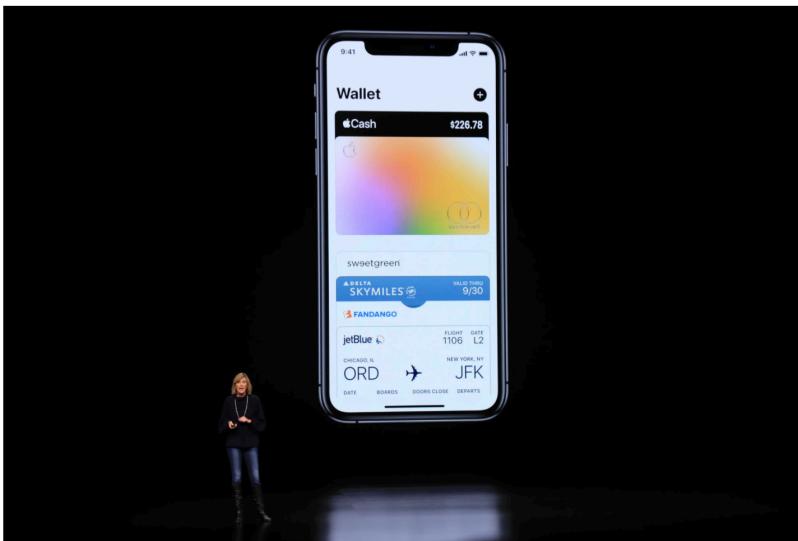
The Chinese government has drawn wide [international condemnation](#) for its harsh crackdown on ethnic Muslims in its western region, including holding as many as a million of them in detention camps.

Now, documents and interviews show that the authorities are also using a vast, secret system of advanced facial recognition technology to track and control the Uighurs, a [largely Muslim minority](#). It is the first known example of a government intentionally using artificial intelligence for racial profiling, experts said.

The facial recognition technology, which is integrated into China's rapidly expanding networks of surveillance cameras, looks exclusively for Uighurs based on their appearance and keeps records of their comings and goings for search and review. The practice makes China a pioneer in applying next-generation technology to watch its people, potentially ushering in a new era of automated racism.

Apple Card Investigated After Gender Discrimination Complaints

A prominent software developer said on Twitter that the credit card was “sexist” against women applying for credit.



Jennifer Bailey, vice president of Apple Pay. Regulators are investigating Apple Card’s algorithm, which is used to determine applicants’ creditworthiness. Jim Wilson/The New York Times

“My wife and I filed joint tax returns, live in a community-property state, and have been married for a long time,” Mr. Hansson wrote Thursday on Twitter. “Yet Apple’s black box algorithm thinks I deserve 20x the credit limit she does.”

Mr. Hansson’s tweets caught the attention of more than just his 350,000 followers.

They [struck a nerve with New York State regulators](#), who announced on Saturday that they would investigate the algorithm used by Apple Card to determine the creditworthiness of applicants.

THE NEW YORK
TIMES
November 2019

DHH  @dhh · Nov 8, 2019

Replying to @dhh

So nobody understands THE ALGORITHM. Nobody has the power to examine or check THE ALGORITHM. Yet everyone we’ve talked to from both Apple and GS are SO SURE that THE ALGORITHM isn’t biased and discriminating in any way. That’s some grade-A management of cognitive dissonance.

DHH  @dhh

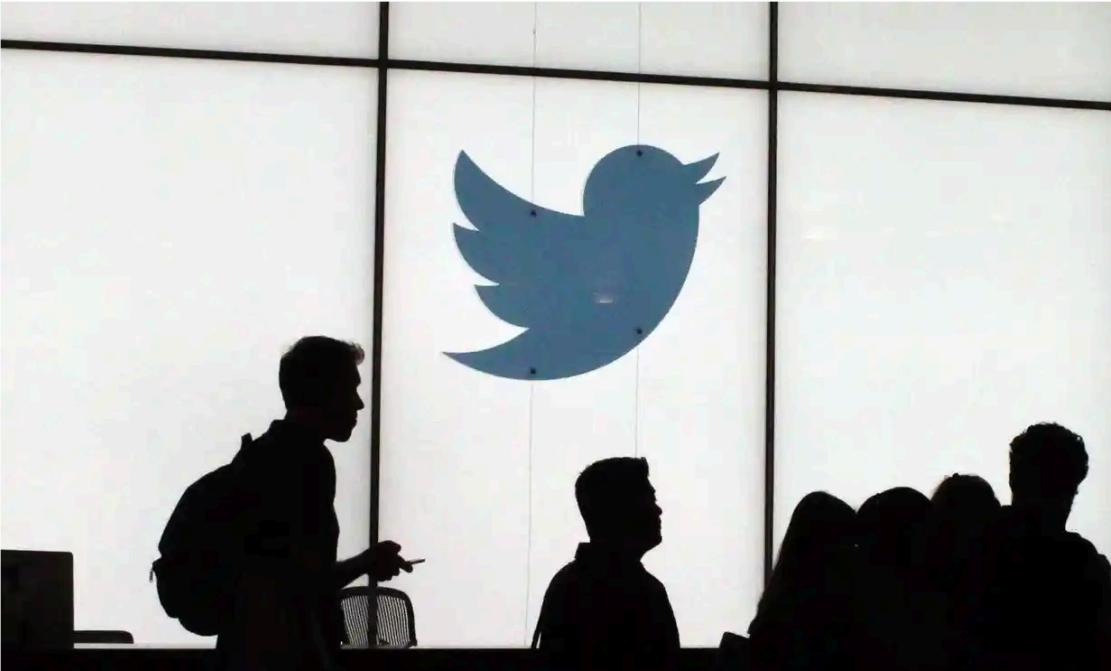
Apple has handed the customer experience and their reputation as an inclusive organization over to a biased, sexist algorithm it does not understand, cannot reason with, and is unable to control. When a trillion-dollar company simply accepts the algorithmic overlord like this...

4:29 PM · Nov 8, 2019 

 4.3K  576 people are Tweeting about this

Twitter apologises for 'racist' image-cropping algorithm

Users highlight examples of feature automatically focusing on white faces over black ones



▲ Twitter users began to spot flaws in the feature over the weekend. Photograph: Glenn Chapman/AFP/Getty Images

[Twitter](#) has apologised for a “racist” image cropping algorithm, after users discovered the feature was automatically focusing on white faces over black ones.



Machine Learning
has progressed
dramatically over
the past decades

However...

Deploying Machine Learning, you face the risk that it be discriminatory, biased, inequitable, exploitative, or opaque.

*We need to figure out how to train models that people can use **safely**, **reliably**, and in a way that they **understand**.*

THE IMPORTANCE OF
INTERPRETABILITY

*Why not just **trust** the model
and **ignore** why it made a
certain decision?*

*“The problem is that **a single metric**,
such as **classification accuracy**,
is an incomplete description of most
real-world tasks.”*

(Doshi-Velez and Kim 2017)

*The necessity for interpretability comes from an **incompleteness** in the problem formalization.*

*For certain problems or tasks it is **not enough** to get the **answer** (the **what**)*

(Doshi-Velez and Kim 2017)

THE CURRENT PROJECT AIMS TO CREATE AN
INTERPRETABLE MACHINE LEARNING MODEL
FOR GENDERrecognition.

SIMILAR RESEARCH

Towards Interpretable Face Recognition

Bangjie Yin^{1*} Luan Tran^{1*} Haoxiang Li^{2†} Xiaohui Shen^{3†} Xiaoming Liu¹

¹Michigan State University ²Wormpex AI Research ³ByteDance AI Lab

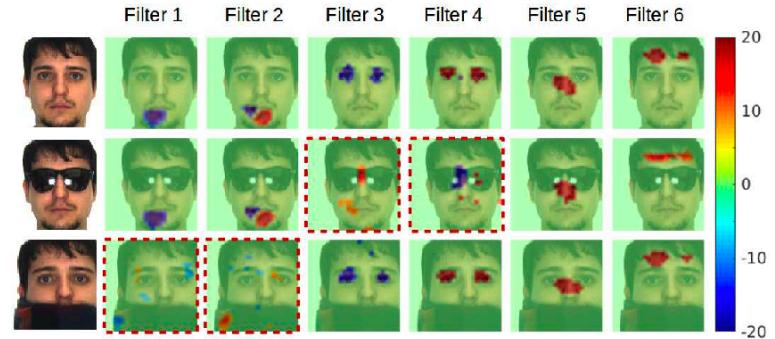
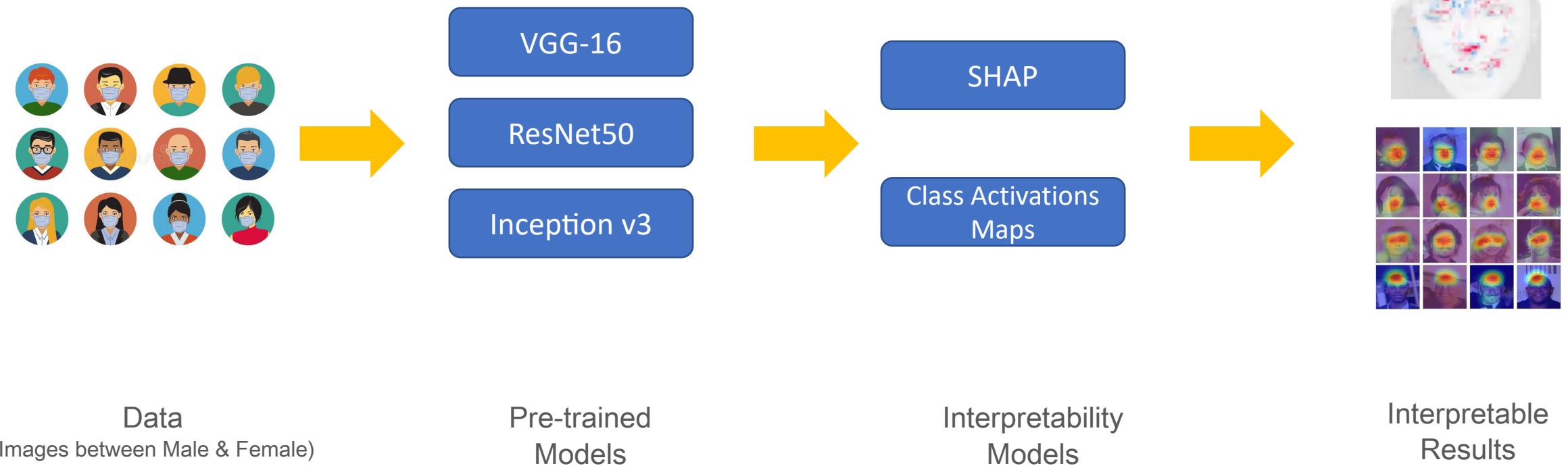


Figure 1. An example on the behaviors of an interpretable face recognition system: left most column is three faces of the same identity and right six columns are filter responses from six filters; each filter captures a clear and consistent semantic face part, e.g., eyes, nose, and jaw; heavy occlusions, eyeglass or scarf, alternate responses of corresponding filters and make the responses being more scattered, as shown in red bounding boxes.

Project Diagram



THE DATA

CelebA Dataset

200k+ face images with attributes
Including gender

84,439
Female: 118,165
Male: 84,439

67,550

Train 80%

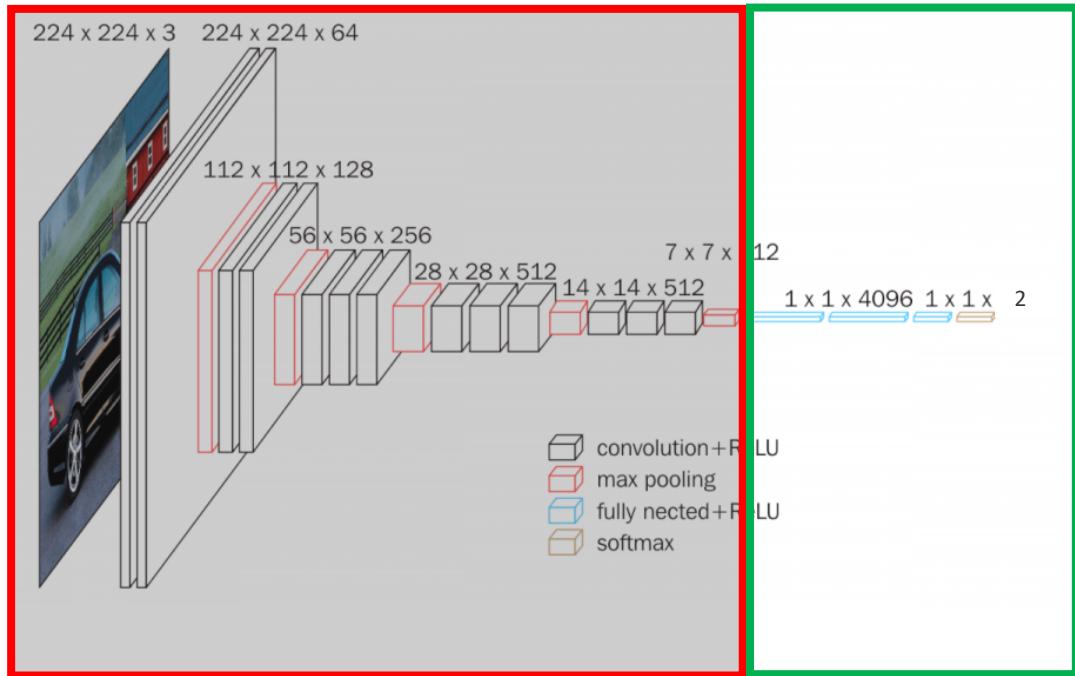
16,889

Test 20%



Pre-trained
MODELS

VGG-16



Source: Inception Module (source: original paper)

Frozen

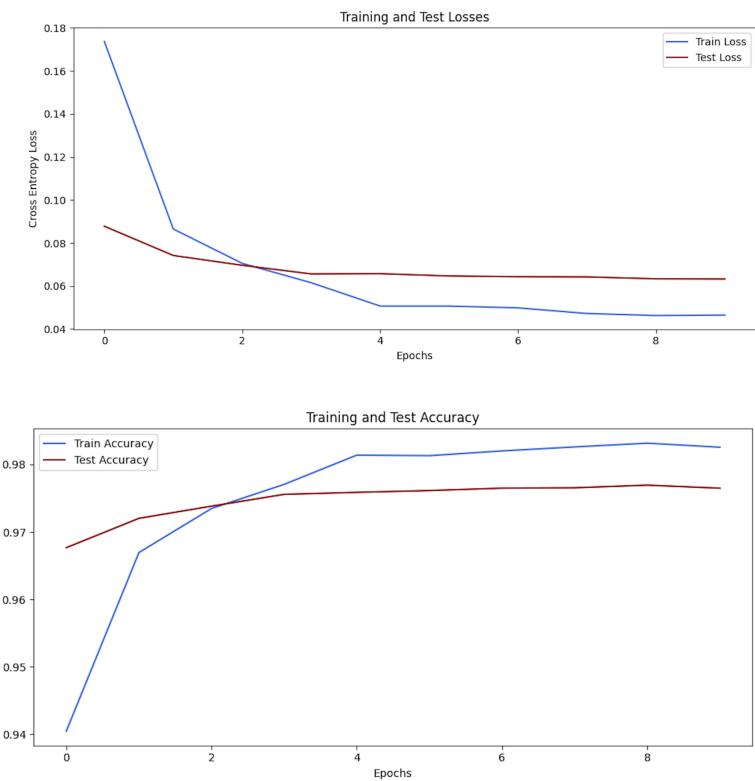
Trainable

Hyperparameters:

Batch size = 64

LR = 0.001

Epochs = 10



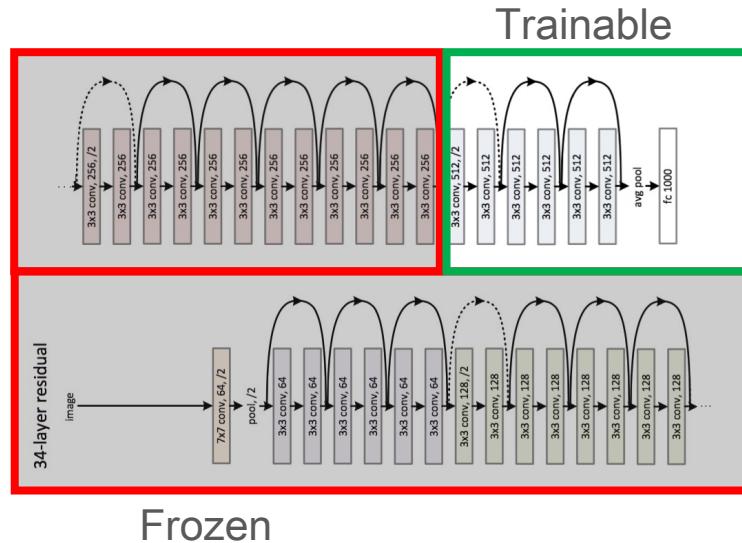
Train Accuracy: 97.6%

Val. Accuracy: 98.2%

Train Loss: 0.05

Val Loss: 0.07

ResNet50



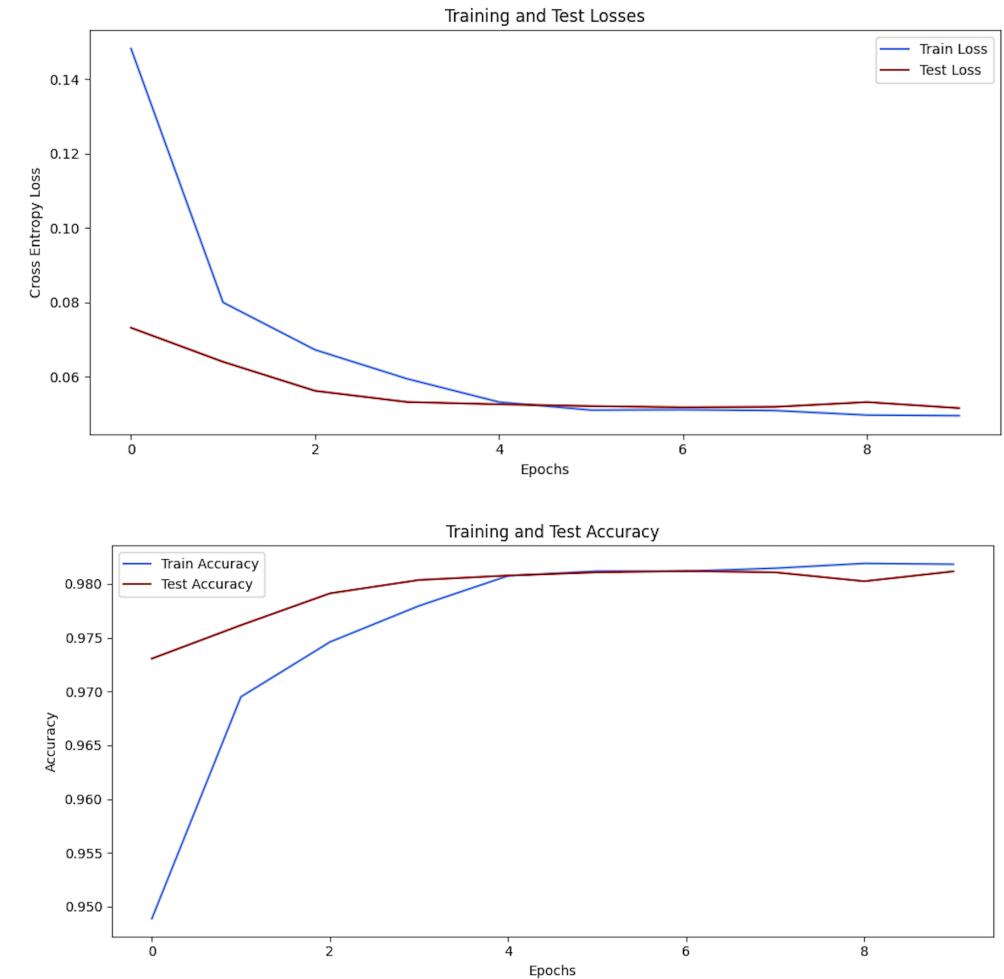
Source: Inception Module (source: original paper)

Hyperparameters:

Batch size = 64

LR = 0.001

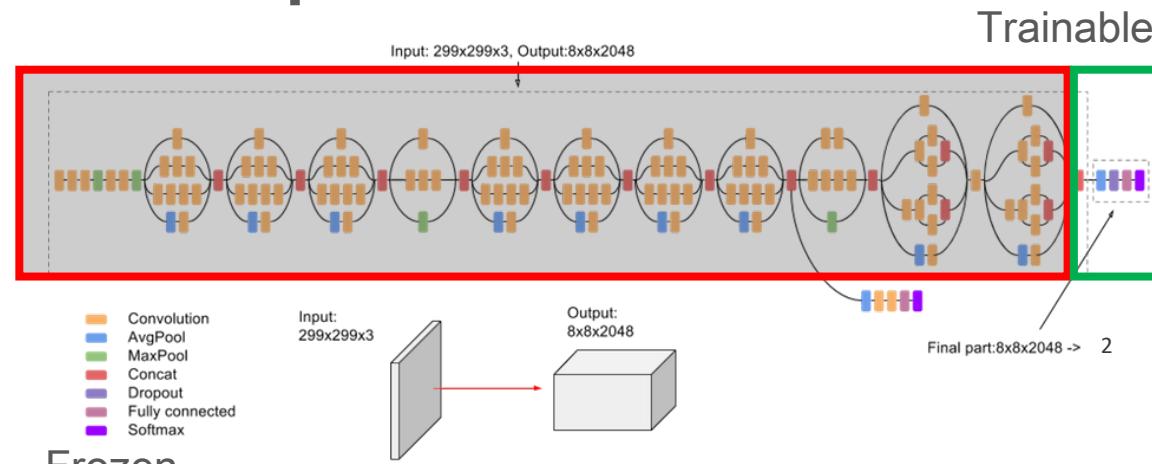
Epochs = 10



Train Accuracy: 98.1 %
Val. Accuracy: 97.9%

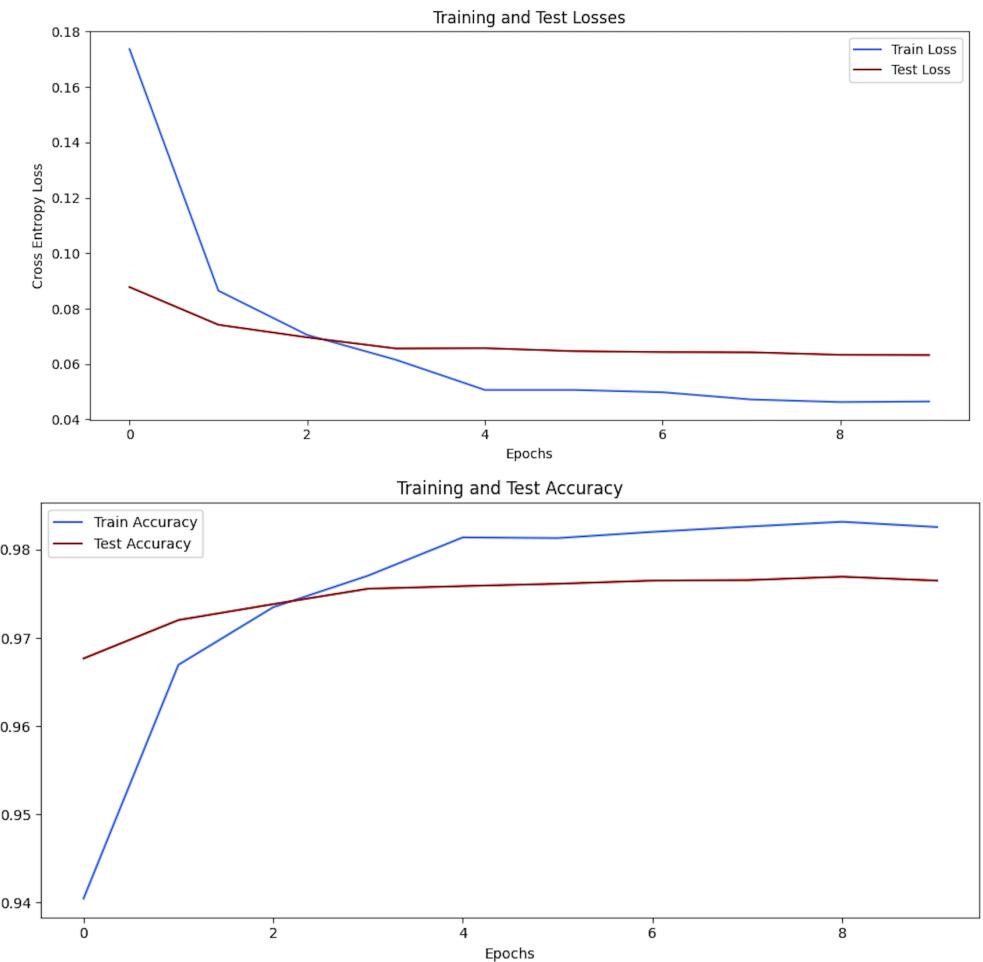
Train Loss: 0.058
Val Loss: 0.054

Inception v3



Source: Inception Module (source: original paper)

Hyperparameters:
Batch size = 64
LR = 0.001
Epochs = 10



Train Accuracy: 98.3%
Val. Accuracy: 97.5%

Train Loss: 0.051
Val Loss: 0.072

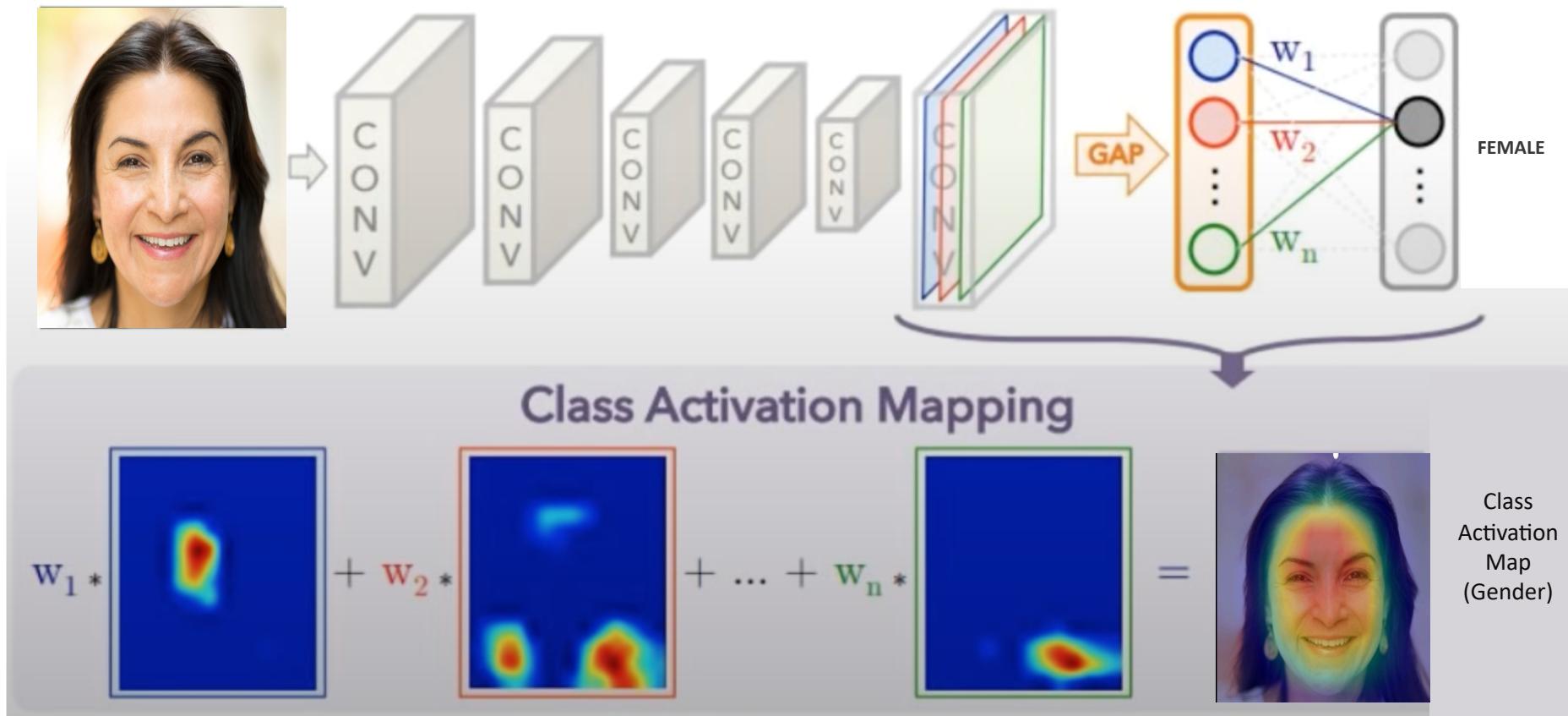
INTERPRETABILITY MODELS

Class Activation Maps

(CAM)

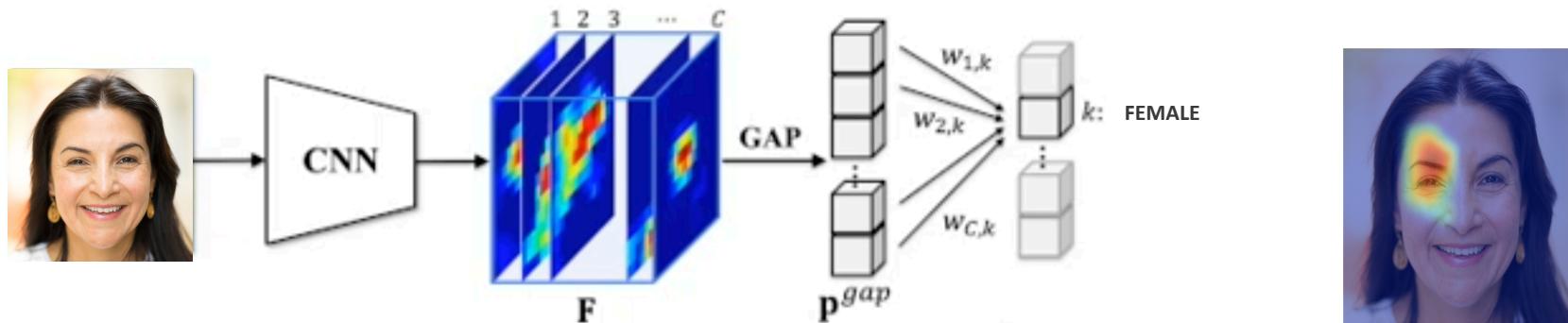
CLASS ACTIVATION MAPS

- The image shrinks, the number of features increases.
- ReLU: All features are positive or zero.



CLASS ACTIVATION MAPS

- We consider one class at a time (the predicted class)
- Example: $W = w[:, \text{human_face_index}]$ # Size 512
- $\text{CAM} = F[0] * w[0] + F[1] * w[1] + \dots + F[512] * w[512]$ (Weighted sum of images)
- Rescale the output image to the original size (i.e. 224×224 for ResNet), and the 2 images over each other



CAM RESULTS FOR FEMALE GENDER

VGG16 focus more on specific features. i.e. eyes, Chin

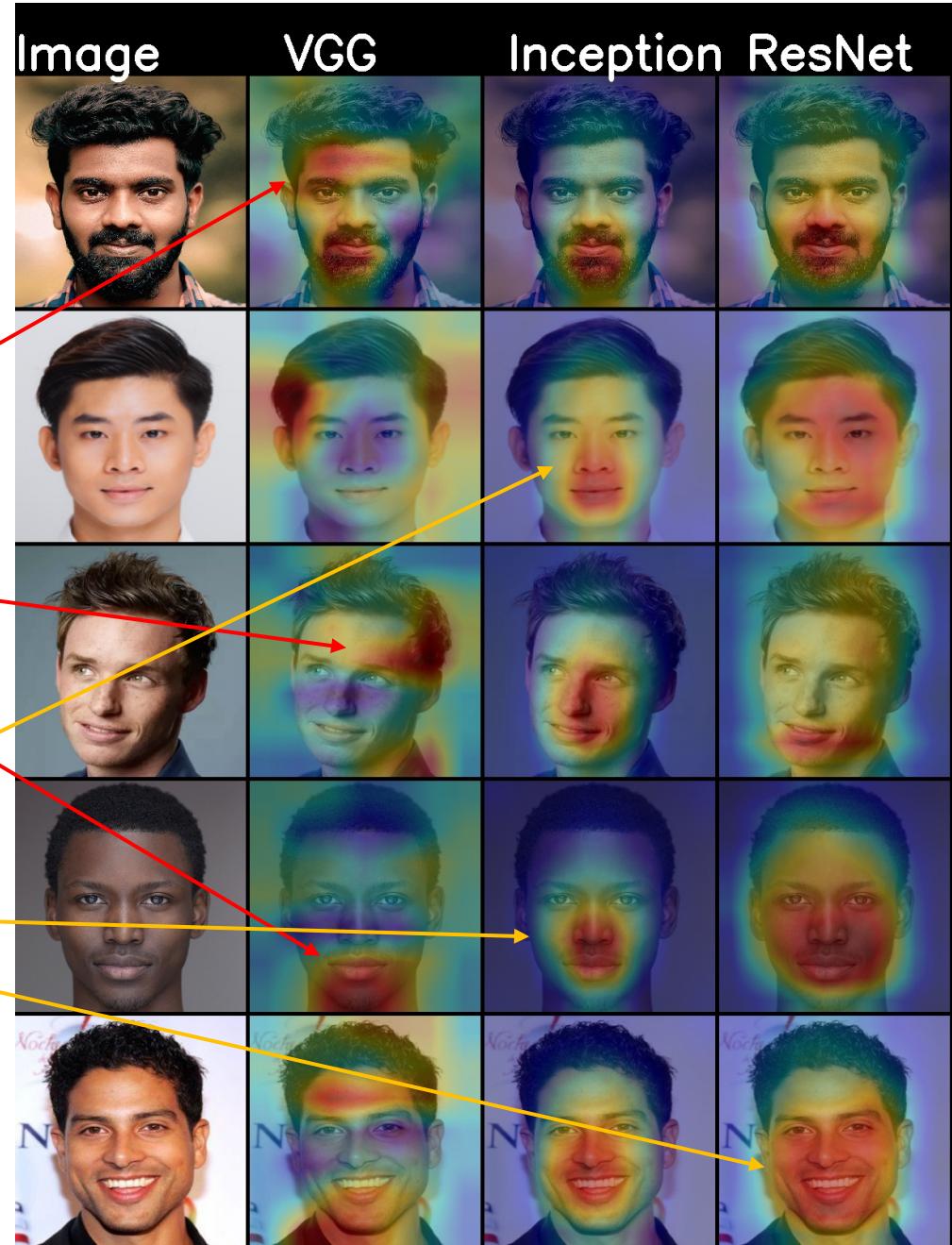
Inception and Resnet Focus on the overall face

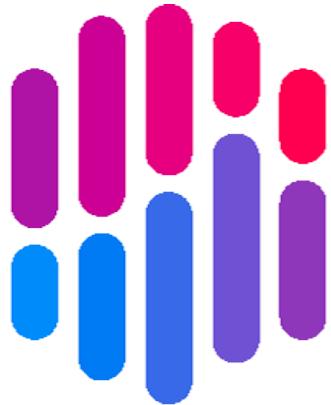


CAM RESULTS FOR MALE GENDER

VGG16 focus more on specific features. i.e. eyes, chin.

Inception and Resnet Focus on the overall face





SHAP

SHAP (SHapley Additive
exPlanations) is a method to explain
individual predictions.

SHAP is based on the game
theoretically optimal Shapley Values.

Lundberg and Lee (2016)

$$g(z') = \phi_0 + \sum_{j=1}^M \phi_j z'_j$$

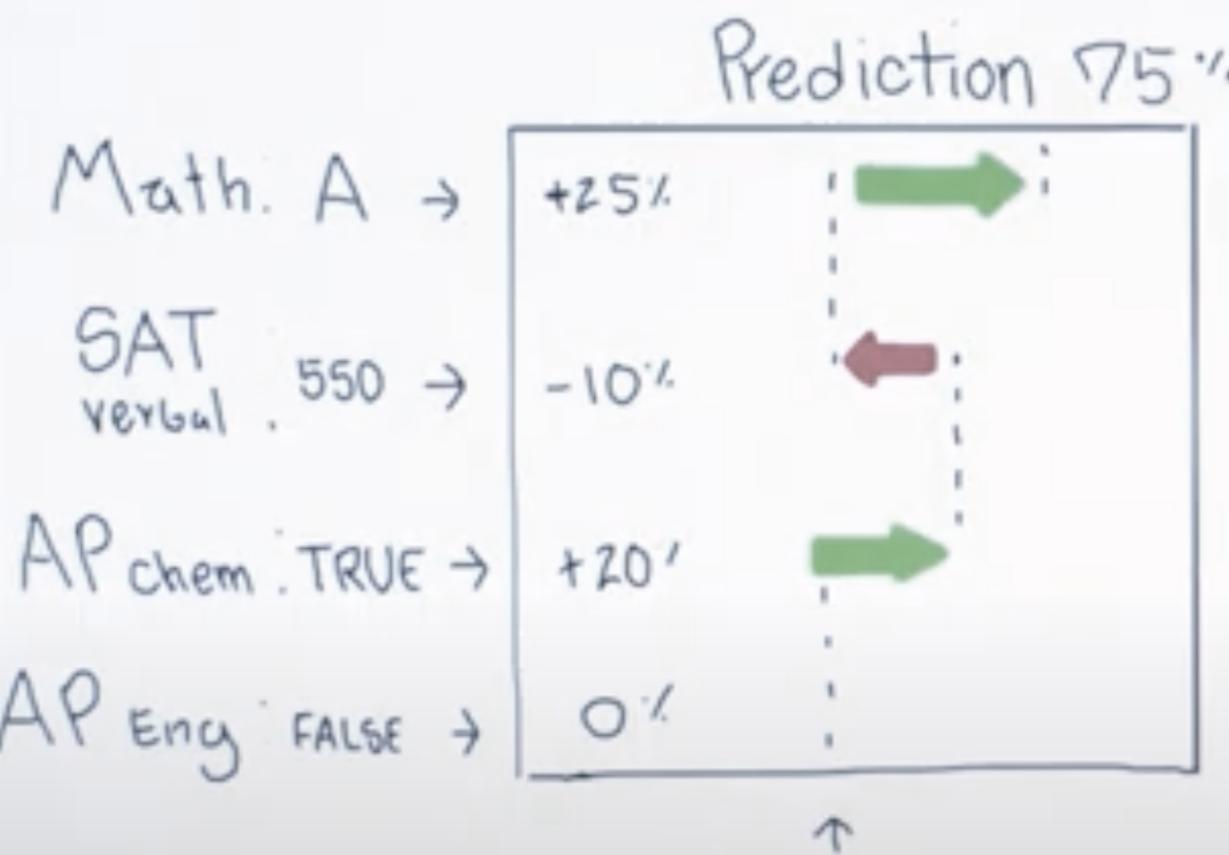
SHAP

- SHAP is based on Game Theory.
- The **game** reproduces the outcome of the model.
- The **players** are the features in the model.
- One game → One observation.
- SHAP quantifies the contribution that each feature brings to the prediction made by the model.

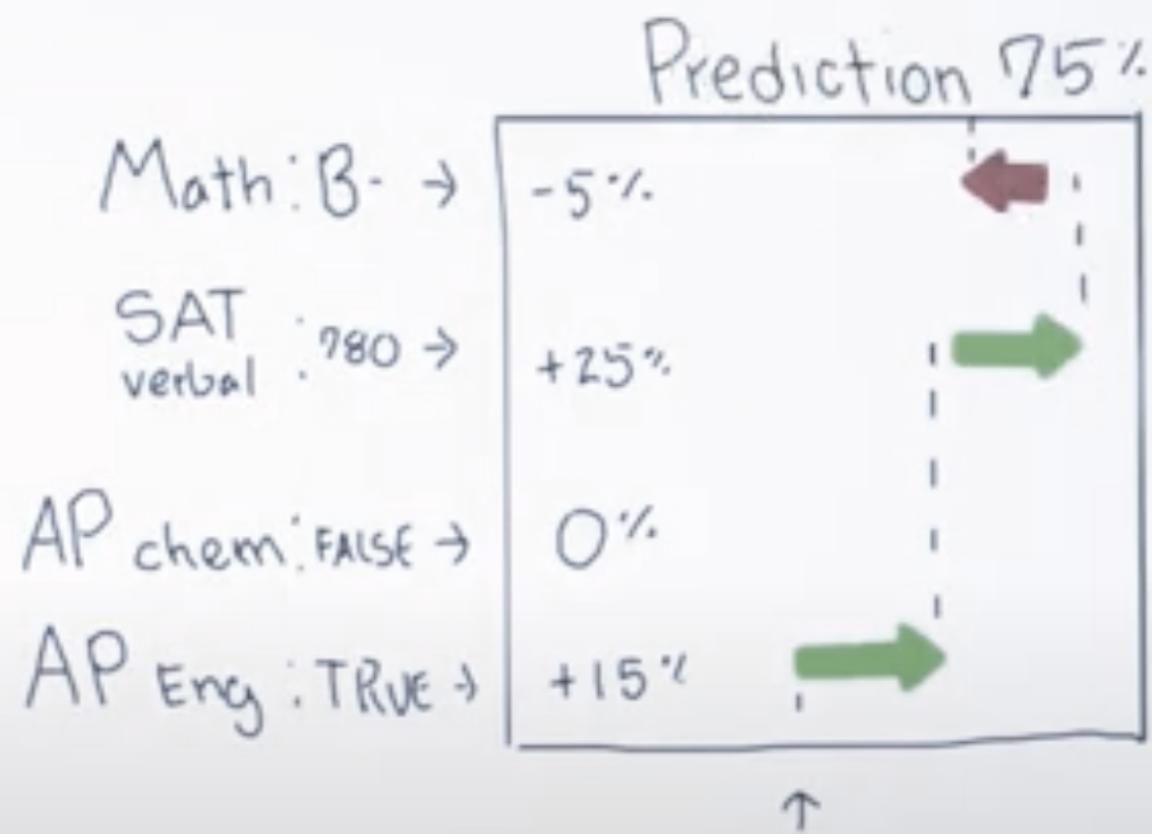
Adam



Ben



Base Rate - 40%



Base Rate - 40%

SHAP

THIS PROJECT WILL USE THE GRADIENT-BASED EXPLANATION METHOD

$$\text{IG}_i(x, x') := (x_i - x'_i) \times \int_{\alpha=0}^1 \frac{\partial F(x' + \alpha \times (x - x'))}{\partial x_i} d\alpha$$

(here $\frac{\partial F(x)}{\partial x_i}$ is the gradient of F along the i^{th} Dimension at x).

- Integrated Gradient (IG) approximate the Shapley values.
- IG is faster Shapley-value-based method.
- Computes the gradients of the model output on a few different inputs (typically 50).

SHAP RESULTS

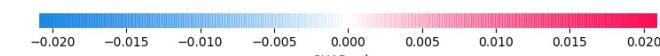
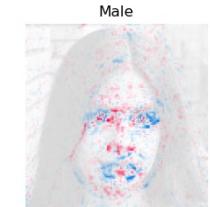
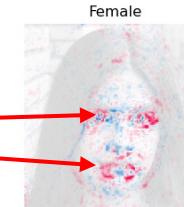
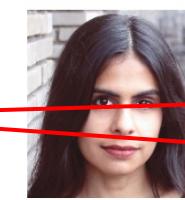
LOW SHAP VALUES HIGH SHAP VALUES

VGG-16

MALE

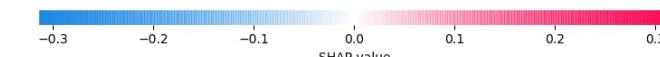
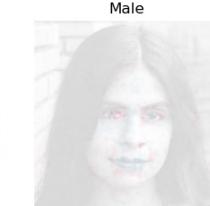
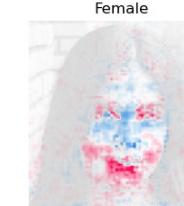
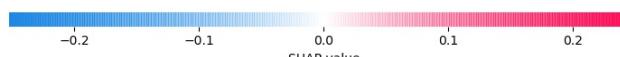
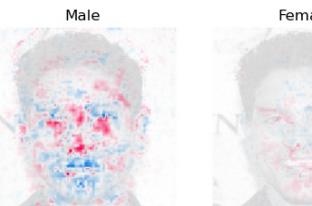


FEMALE

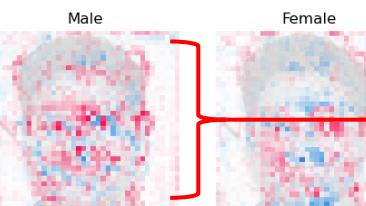


Focus on specific features

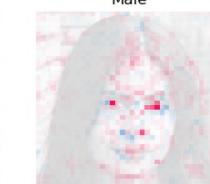
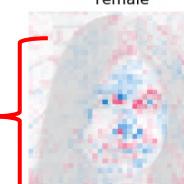
INCEPTION v3



RESNET



Focus on overall image



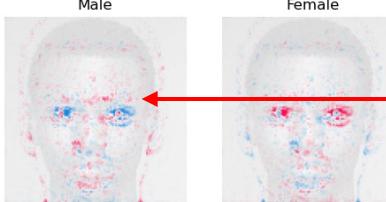
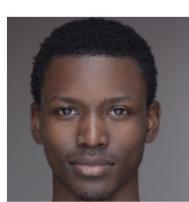
SHAP RESULTS

LOW SHAP VALUES

HIGH SHAP VALUES

VGG-16

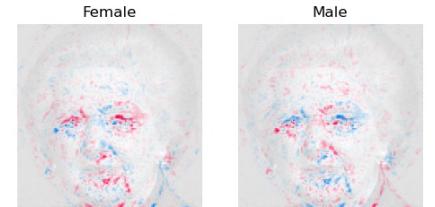
MALE



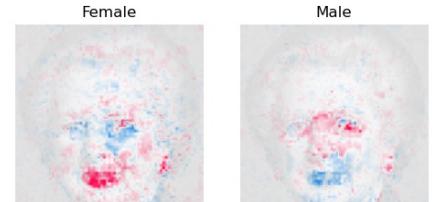
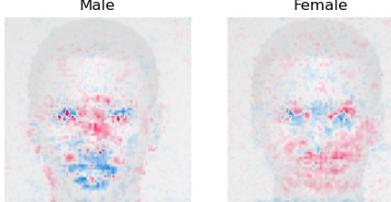
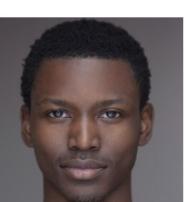
VGG16 is showing
low SHAP (blue)
values instead of high
values (red)



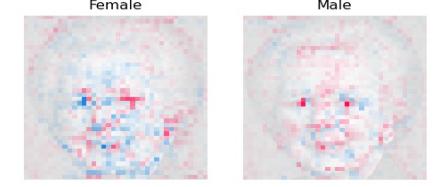
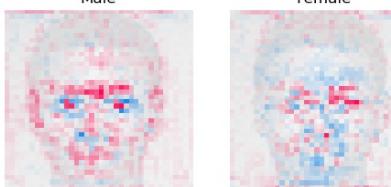
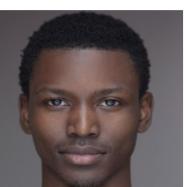
FEMALE



INCEPTION

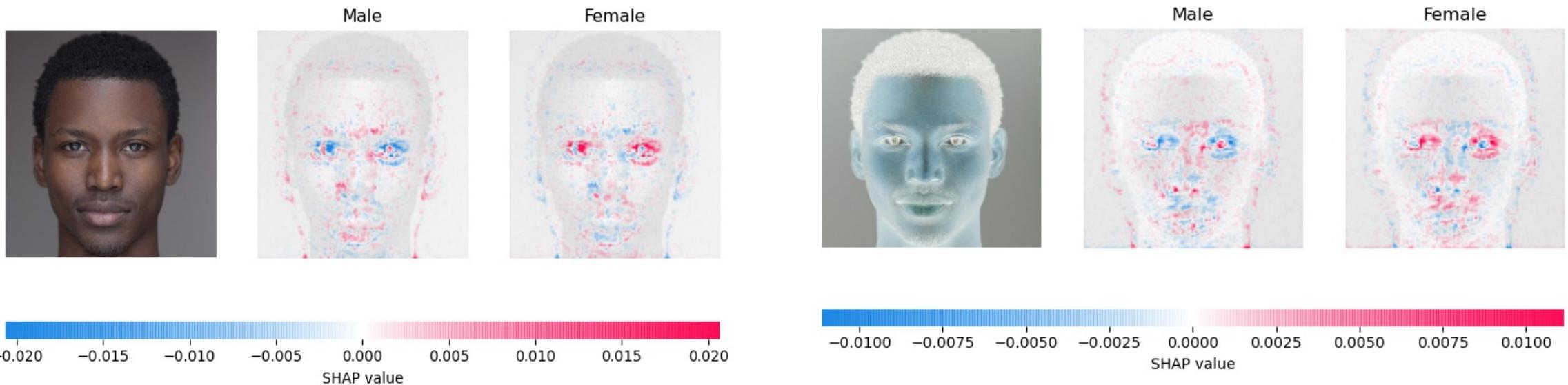


RESNET



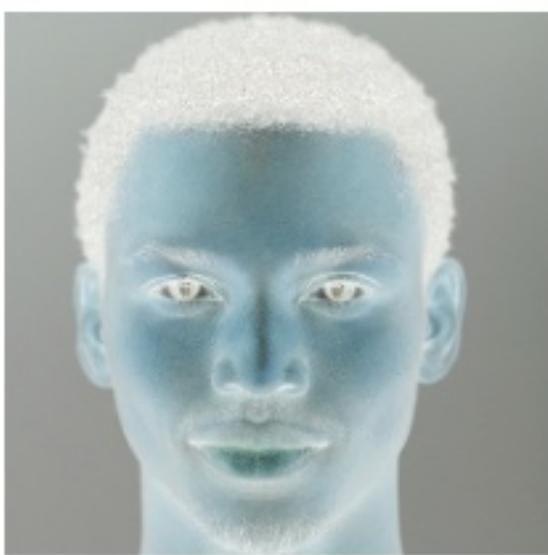
SHAP RESULTS

LOW SHAP VALUES HIGH SHAP VALUES



SHAP RESULTS

LOW SHAP VALUES HIGH SHAP VALUES



Male



Female



COMPARING CAM VS SHAP

(Gender – Female)



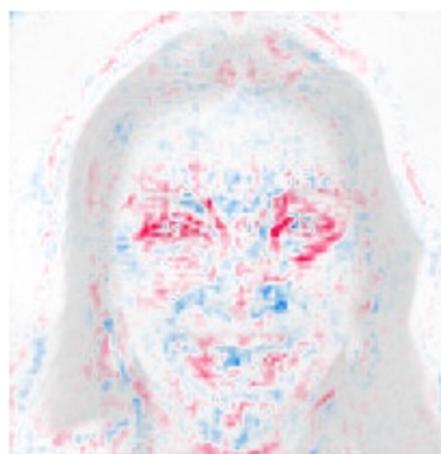
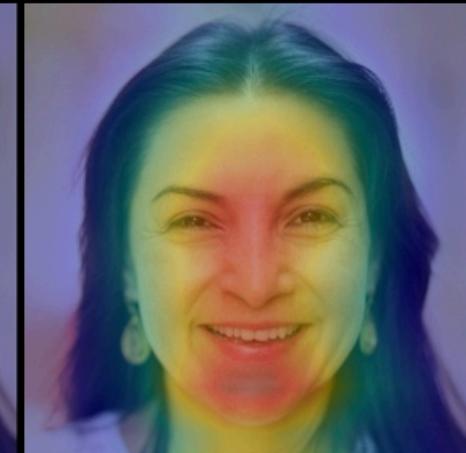
VGG-16



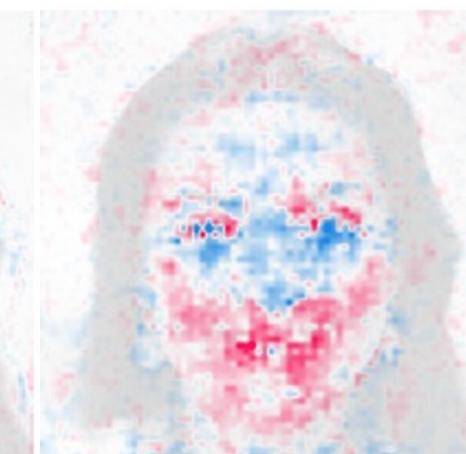
Inception v3



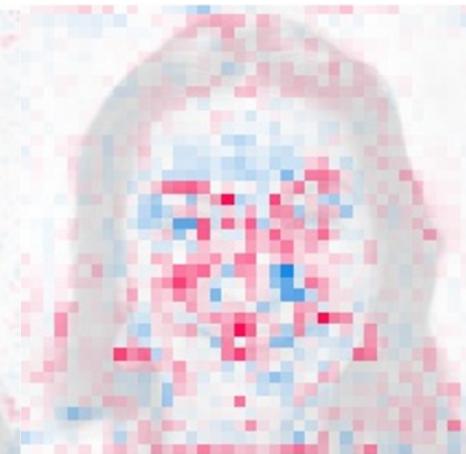
ResNet



-0.015 -0.010 -0.005 0.000 0.005 0.010 0.015
SHAP value



-0.15 -0.10 -0.05 0.00 0.05 0.10 0.15
SHAP value



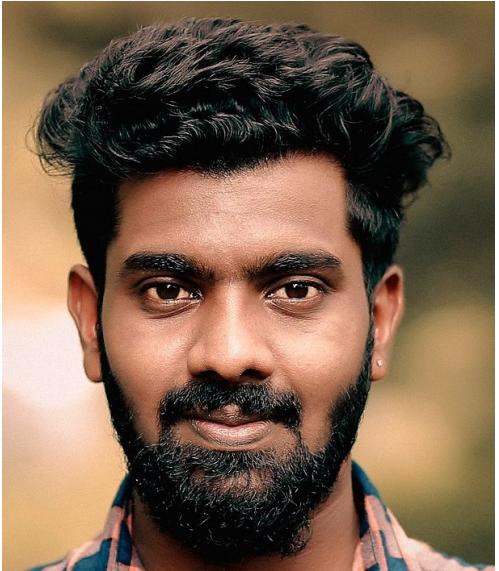
-0.15 -0.10 -0.05 0.00 0.05 0.10 0.15
SHAP value

LOW SHAP VALUES

HIGH SHAP VALUES

COMPARING CAM VS SHAP

(Gender – Male)



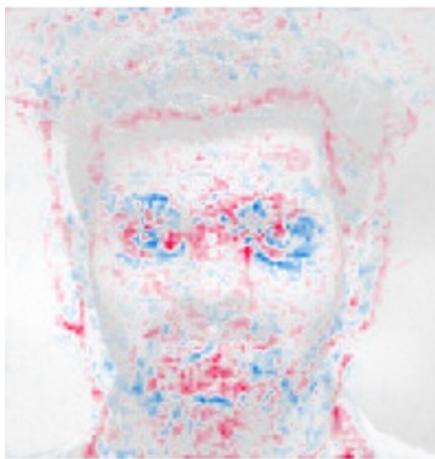
VGG-16



Inception



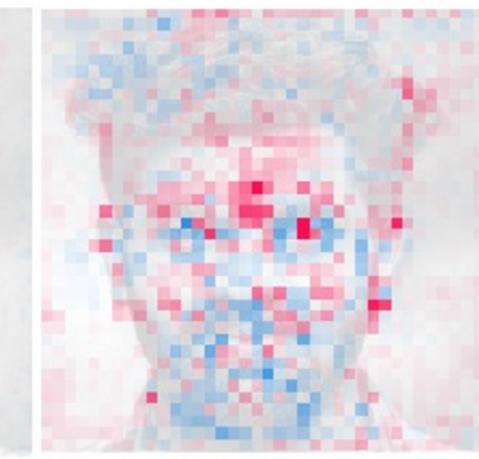
ResNet



-0.015 -0.010 -0.005 0.000 0.005 0.010 0.015
SHAP value



-0.15 -0.10 -0.05 0.00 0.05 0.10 0.15
SHAP value



-0.10 -0.05 0.00 0.05 0.10
SHAP value



LOW SHAP VALUES



HIGH SHAP VALUES

COMPARING CAM VS SHAP

(Gender – Female)



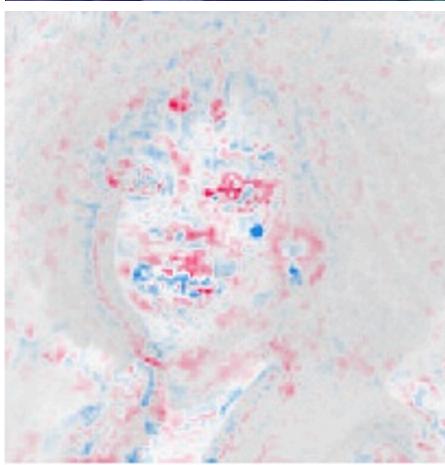
VGG-16



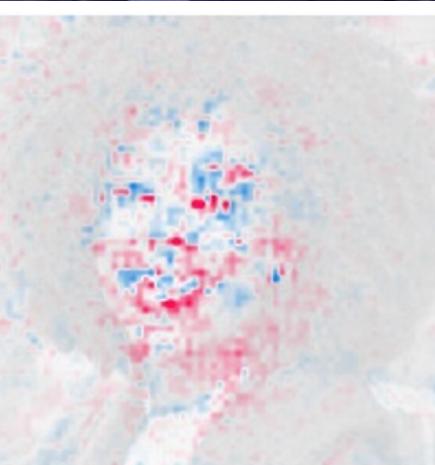
Inception



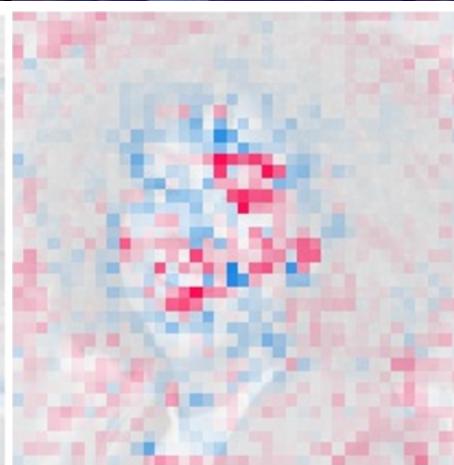
ResNet



-0.015 -0.010 -0.005 0.000 0.005 0.010 0.015
SHAP value



-0.15 -0.10 -0.05 0.00 0.05 0.10 0.15
SHAP value



-0.10 -0.05 0.00 0.05 0.10
SHAP value



LOW SHAP VALUES



HIGH SHAP VALUES

COMPARING CAM VS SHAP

(wearing sunglasses)



VGG-16



Inception



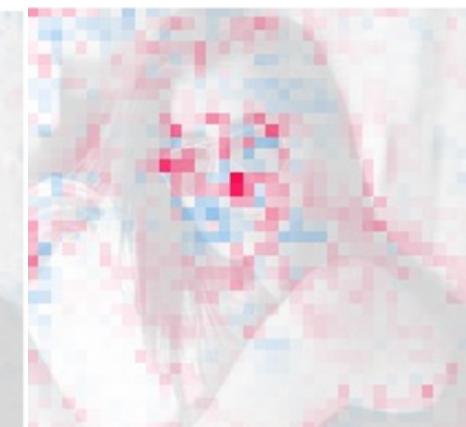
ResNet



-0.02 -0.01 0.00 0.01 0.02
SHAP value



-0.15 -0.10 -0.05 0.00 0.05 0.10 0.15
SHAP value



-0.15 -0.10 -0.05 0.00 0.05 0.10 0.15
SHAP value



LOW SHAP VALUES



HIGH SHAP VALUES

COMPARING CAM VS SHAP

(smiling)



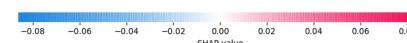
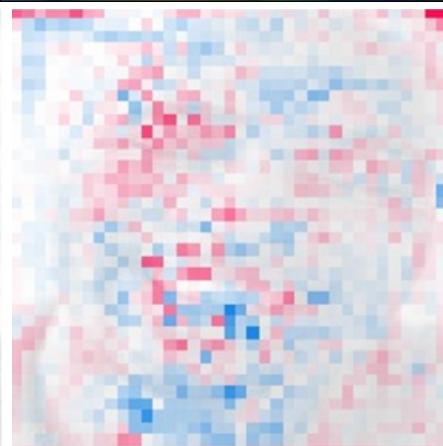
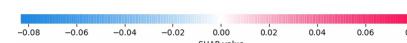
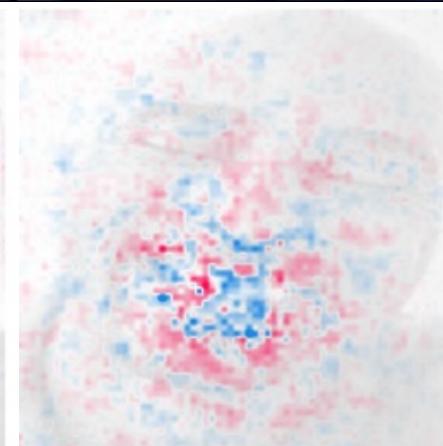
VGG-16



Inception



ResNet



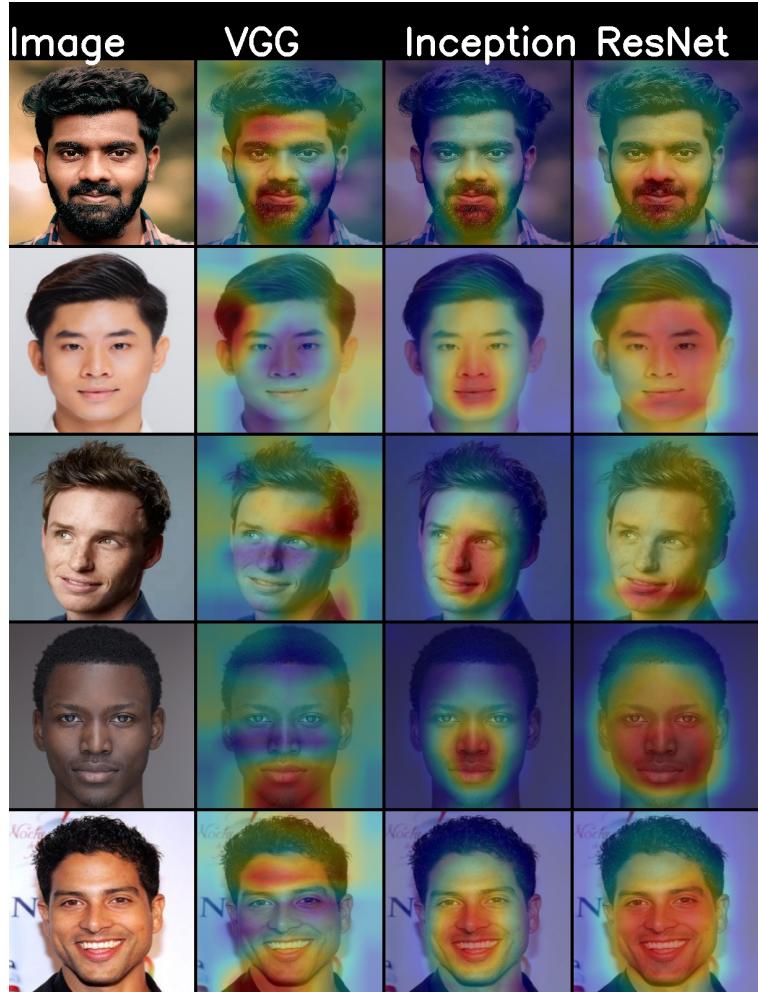
LOW SHAP VALUES



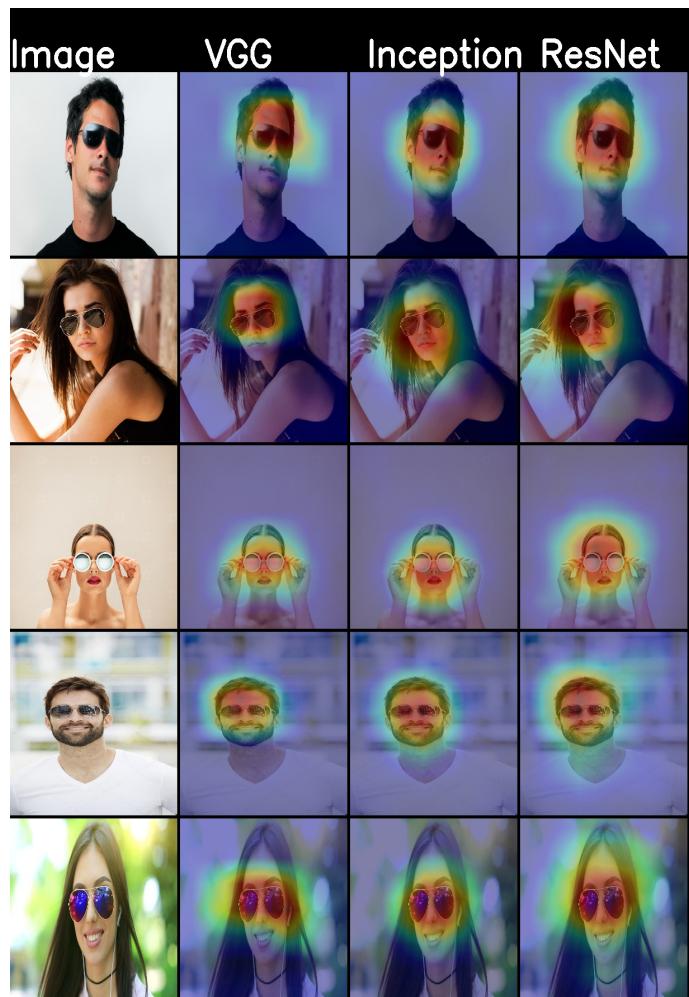
HIGH SHAP VALUES

IDENTIFYING BIAS

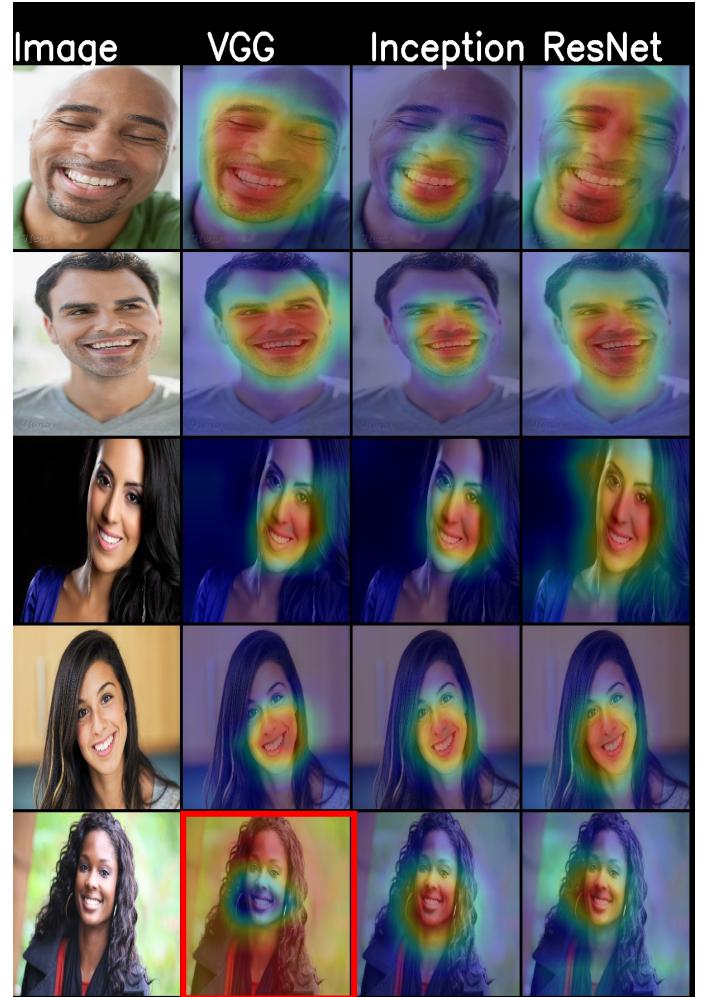
GENDER MALE



WEARING SUNGLASSES

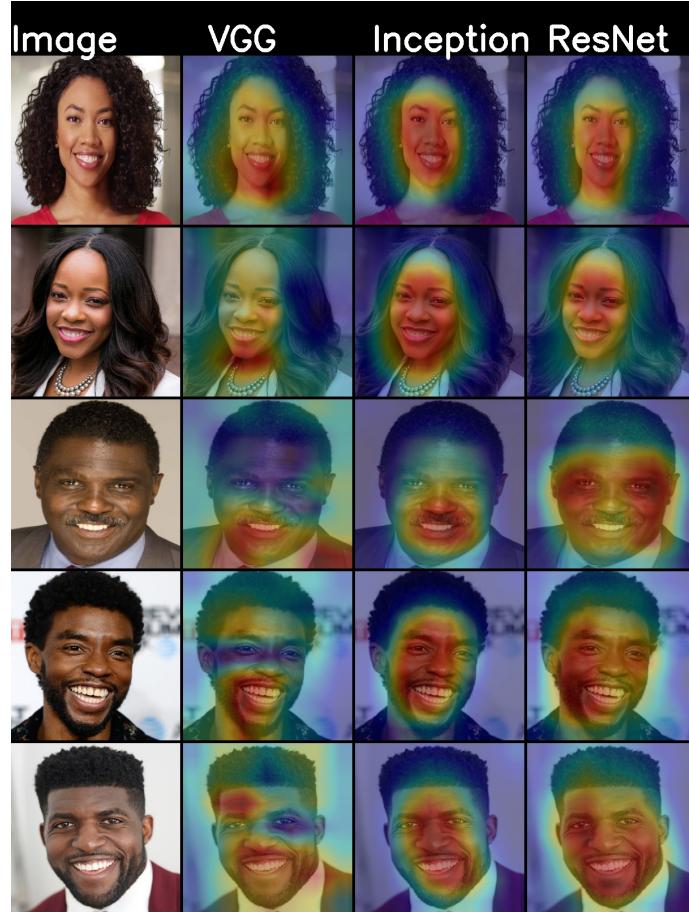
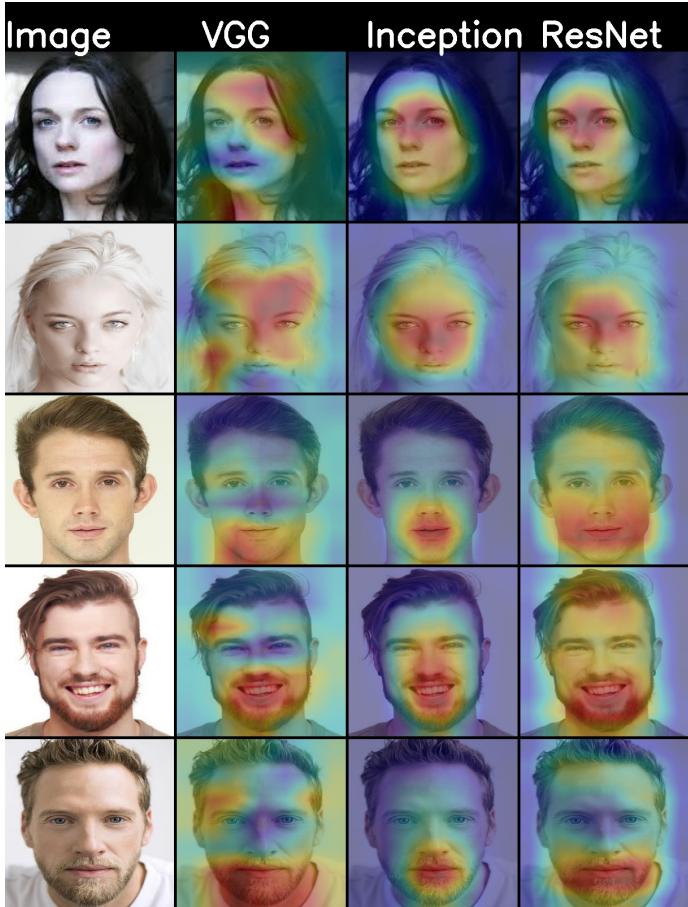


SMILING



The most important feature,
turns out to be the less important

IDENTIFYING BIAS



HARD TO CONCLUDE THAT THE MODELS ARE BIASED

CONCLUSIONS

- CAM can help in visualizing the predicted class, highlighting the features detected by the CNN.
- Interpretability can give certainty to complex models to make correct predictions in an ethical way.
- It can help to spot bias in the model.
- Can Help you explain your model.
- In some countries, regulation can make it mandatory.

REFERENCES

- Zhou, B., Khosla, A., Lapedriza, A., & Oliva, A. (2016). Learning Deep Features for Discriminative Localization. Retrieved from http://cnnlocalization.csail.mit.edu/Zhou_Learning_Deep_Features_CVPR_2016_paper.pdf
- Petsiuk, Vitali, et al. *RISE: Randomized Input Sampling for Explanation of Black-Box Models*, 2016.
- Käärkäinen Kimmo, and Jungseock Joo. *FairFace: Face Attribute Dataset for Balanced Race, Gender, and Age*, Aug. 2019.
- <https://news.mit.edu/2020/aleksander-madry-machine-learning-0308>
- <https://www.pewresearch.org/internet/2019/09/05/the-challenges-of-using-machine-learning-to-identify-gender-in-images/>
- <https://science.sciencemag.org/content/349/6245/255.full>
- https://shap.readthedocs.io/en/latest/example_notebooks/gradient_explainer/Multi-input%20Gradient%20Explainer%20MNIST%20Example.html
- <https://blog.fiddler.ai/2019/08/should-you-explain-your-predictions-with-shap-or-ig/>
- <https://towardsdatascience.com/shap-explained-the-way-i-wish-someone-explained-it-to-me-ab81cc69ef30>
- <http://cvlab.cse.msu.edu/project-interpret-FR>