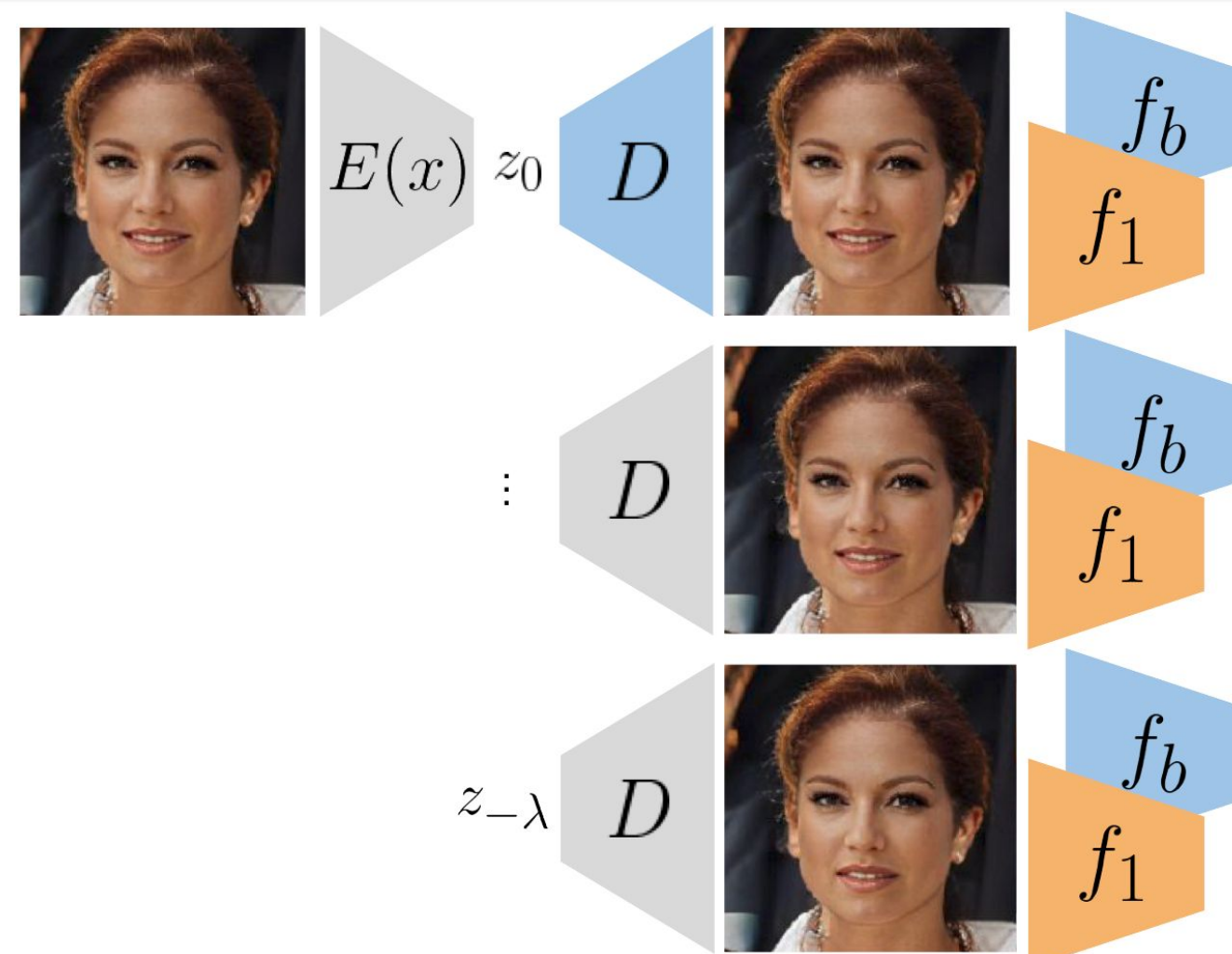# Identifying Spurious Correlations using Counterfactual Alignment
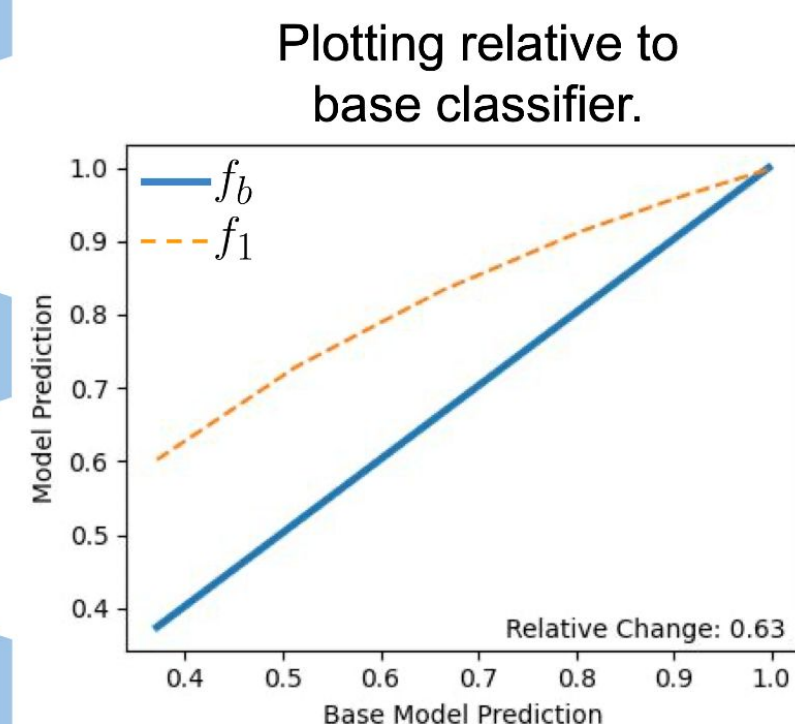
Joseph Paul Cohen, Louis Blankemeier, Akshay Chaudhari

Stanford

## Counterfactual Alignment



Plotting relative to base classifier.

Images generated at different lambdas w.r.t. the base classifier.

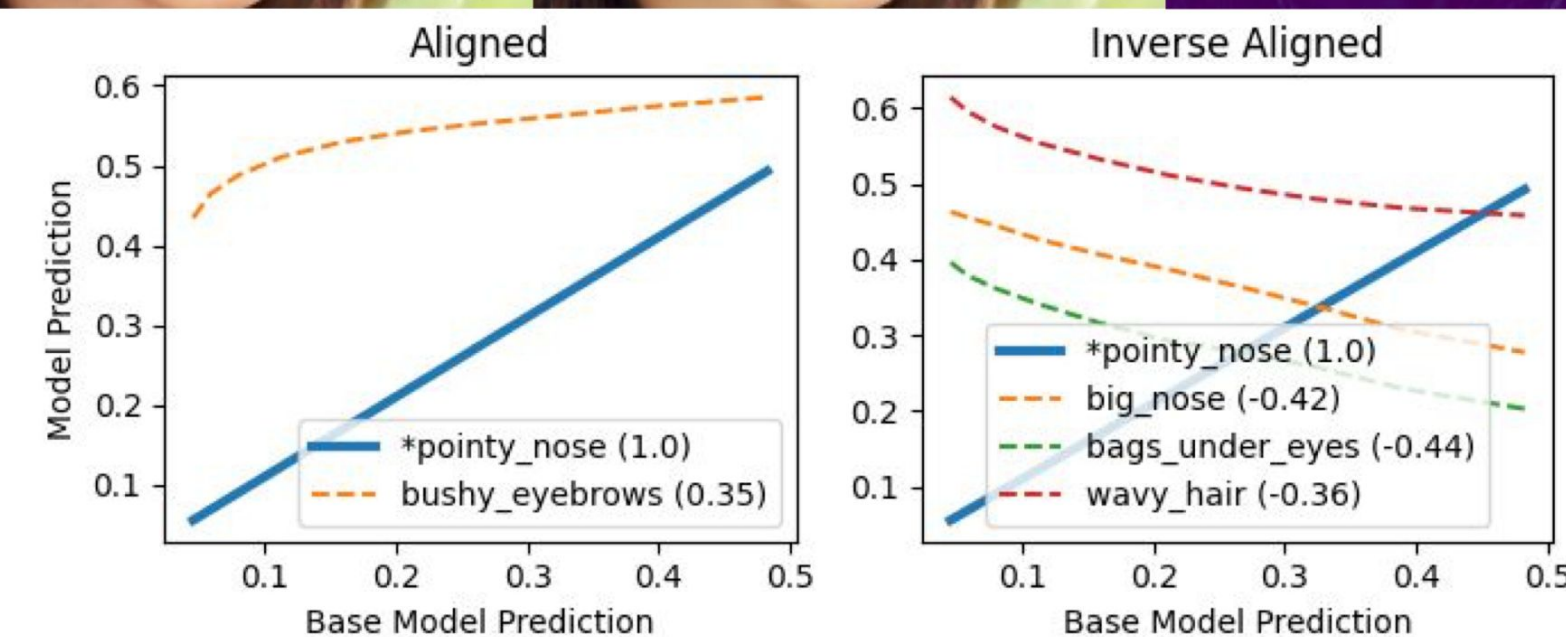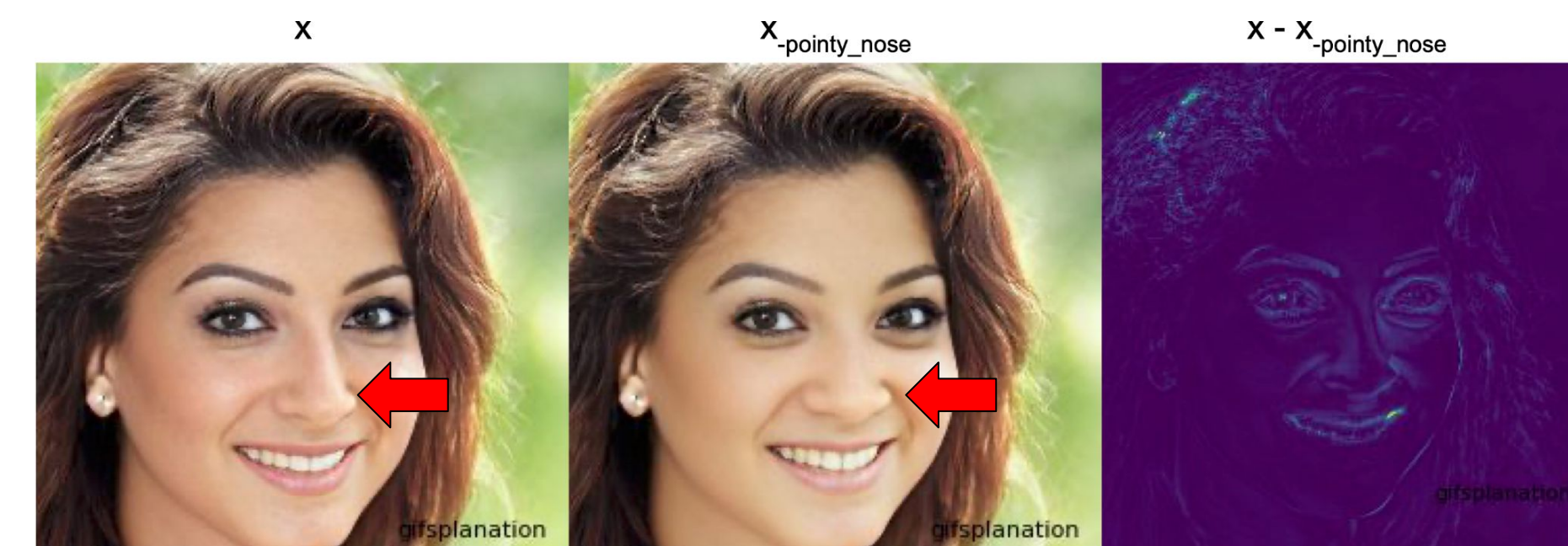Images processed with multiple classifiers.

Compute relative change between the predictions:

$$\frac{f_1(D(z_\lambda)) - f_1(D(z_0))}{f_b(D(z_\lambda)) - f_b(D(z_0))}$$

- Counterfactuals are generated by subtracting the gradient of the classifier output w.r.t. the latent representation.
- The representation is reconstructed back into multiple images with different magnitudes of change (depending on the $\lambda$).
- Reconstructed images are processed with multiple classifiers.
- Base model predictions can be used as the x-axis to more easily compare it to the predictions of another classifier.
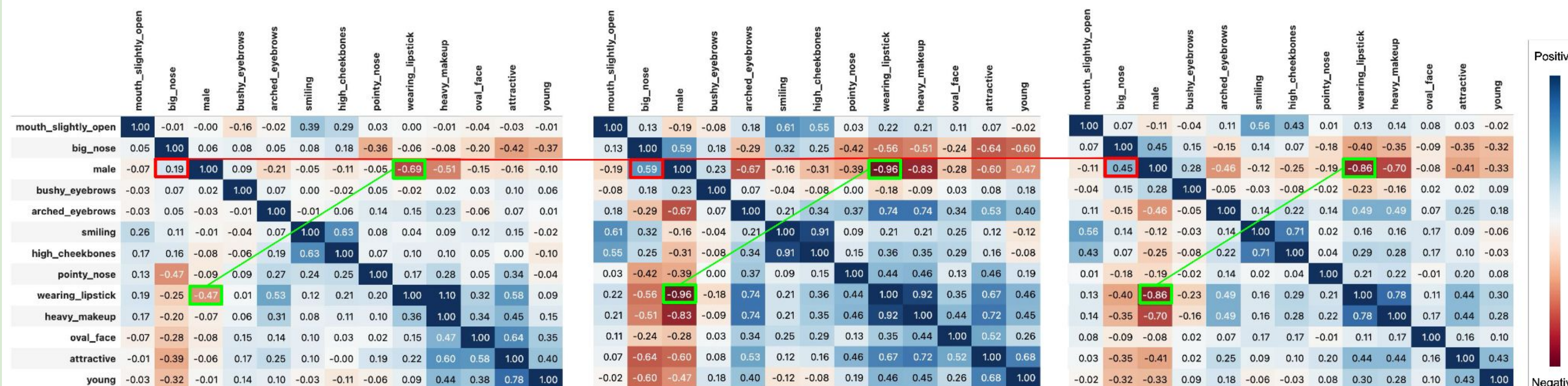- Changes can be quantified and compared using relative change.

## Studying Face Attribute Classifiers



CF alignment for pointy nose showing an inverse alignment with big nose and potential spurious relationships with eyebrows eyes and hair. The relative change is shown next to each classifier name.
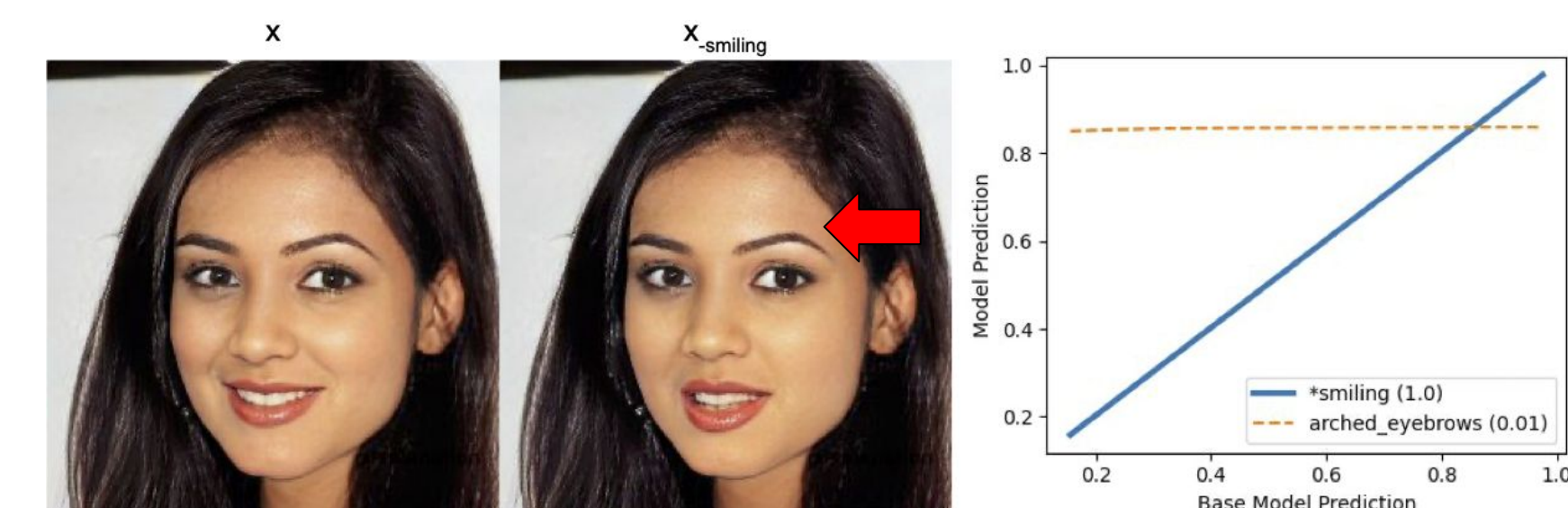
## Aggregate Spurious Correlations

Relationships between face attribute classifiers using different metrics.
Base classifiers are along the rows and downstream classifiers are along the columns.



a) Relative change using CF alignment.
b) Correlation between predictions of classifier.
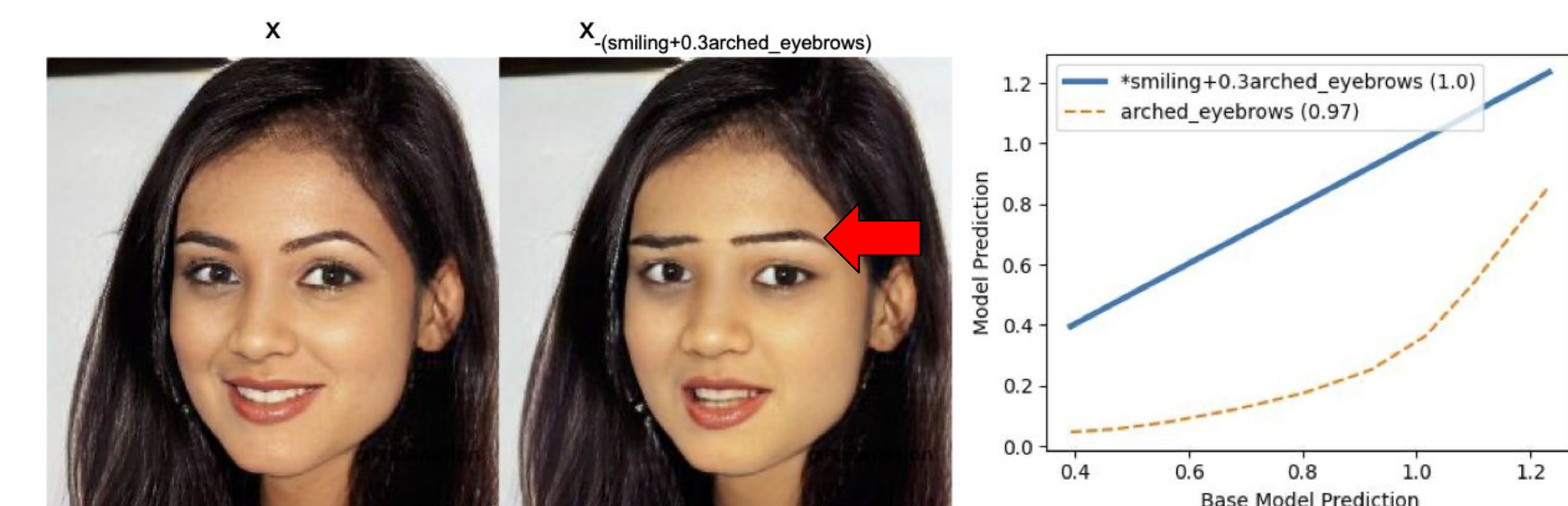c) Correlation between labels in training data.

- Comparing (a) to (b) and (c), shows that many relationships reflected in the CF outputs are preserved from correlations in the training data.
- A non-symmetric relationship between male and wearing lipstick (highlighted in green) reflects an expected relationship.
- The relationship between male and big nose (highlighted in red) is strong in both the classifier predictions and ground truth labels but low in CF alignment, indicating that although correlated, these features aren't exploited by the classifier.

## Detecting Induced Spurious Correlations

Example detecting a spurious correlation in a biased classifier. A classifier is biased with arched eyebrows and this is observed in the CF alignment plot as well as in the counterfactual image. We can observe that in order to decrease the prediction of the (now biased) classifier, the CF image also removes the arched eyebrows.



CF for smiling showing eyebrows are unchanged. The horizontal line indicates the prediction of arched eyebrows is not influenced by the features used for smiling in this image.



CF for the modified smiling classifier with an arched eyebrow spurious correlation showing that eyebrows are now changed.