# CS271P Sponsored Project

Daniel Swartz  Qi Song  Zhongyao Qian

December 2021

**Abstract**

Within the last decade, AI techniques have become particularly useful in assisting with criminal justice needs. One of the most popular AI techniques is facial recognition, which can help establish an individual's identity and trace their relationships through photos. In this project, a CNN model is trained with the Families in the Wild dataset to analyze the paired images of portraits. The resulting models can accurately classify the kinship between two blood-related persons.
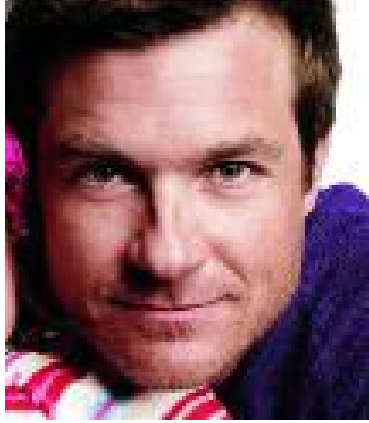
## 1 Introduction

In a wide variety of fields, from genealogy to criminal investigations, it is crucial to determine how people are related to each other. Advancements in computer vision and Convolutional Neural Networks (CNNs) allow us to accelerate the development of creating and documenting relationships between individuals, solely based on images of their faces. Employing AI greatly reduces the intrusiveness of making familial ties during an investigation and the manhours required to infer these connections. We believe that knowing exactly how individuals in a family are related is important, so we implement and fine-tune a CNN model to determine the specific familial ties between portrait images of individuals within the same family.

The following paper is organized into sections. Section 2 covers the dataset and how it is used. Section 3 details how the baseline model and ablation studies are constructed and how training and validation are conducted. Section 4 contains the results and analysis regarding these findings. Lastly, Section 5 concludes the overall project.

The Jupyter Notebook used for this project is saved as "cs271p.ipynb" in the root directory of the Google Drive. drive.google.com/drive/folders/ 19X19TdBt0zFSAhtdIaFuPOIWLbvuG1A-?usp=sharing

## 2 Data

The dataset used is the Families in the Wild [1] image database, which is structured into three parts: training data, validation data, and test data. The training and validation data is structured in folders of families, numbering 1018

(a) Person one.        (b) Person two.

Figure 1: Examples of the types of portraits found within the Families in the Wild dataset.

different groups. Each folder contains sub-folders for each family member, containing one or more cropped images of the face of each individual [Figure 1]. The photos for each member can vary in the age of that person. For example, there can be images from decades ago and from within the last 5 years for a given person. The family folders also contain a comma-separated value (CSV) list of the relationships between each member, each member's name, and their gender. Additionally, there are two separate CSV tables for the training and validation data that contain the information detailing the family number, the two member IDs being matched, what their relationship tie is, and how many sample images there are for each person. The relationships consist of eleven different class labels, bb (brother-brother), fd (father-daughter), fs (father-son), gfgd (grandfather-granddaughter), gfgs (grandfather-grandson), gmgd (grandmother-granddaughter), gmgs (grandmother-grandson), md (mother-daughter), ms (mother-son), ss (sister-sister), and sibs (brother-sister).

Upon further studying of the dataset, we realized that this set is imbalanced in terms of the types of relationship pairs. The percentages of each combination in the training data are illustrated in Figure 2, with the highest percentage being ms at 24.5% of the total relationships and the lowest being gmgd with 0.9%. Additionally, the validation data consists of a smaller dataset but is also imbalanced. The highest percentage in the validation pairs is fs at 18.7%, and the lowest percentages are gmgs and gmgd at 1.4% [Figure 3].

The images for both the training and validation data are preprocessed by loading them as a JPEG image, converting the image to grayscale, and resizing the image to 108 by 124 pixels. This function serves to standardize the images, as the original images come in different sizes and have different hues and lighting. Following the preprocessing, two new datasets, one for training and one for

validation, are created. Both datasets consist of concatenated pairings between an image of the first person and second person that are to be paired, along with the label of their relationship. This creates a two-channel tensor with the shape of (128, 108, 124, 2). The batch size was selected to be 128. The order of the pairs is then shuffled to ensure that the model does not learn the order of the images and overfit to a certain pattern.
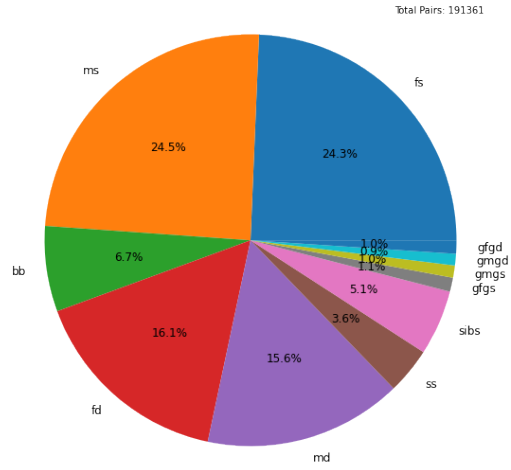
**Percentages of Familial Relationships in Training Data**

Total Pairs: 191361

Figure 2: Chart of the percentages for each pair in the training data.
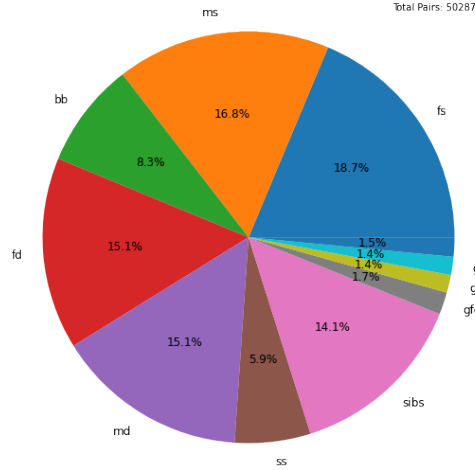
Figure 3: Chart of the percentages for each pair in the validation data.

## 3 Models

### Baseline

For the baseline model, we decided on using a Convolutional Neural Network (ConvNet) with an EfficientNet architecture, specifically EfficientNetB0. The EfficientNet family boasts both higher accuracy and efficiency than most other ConvNets [2]. We perform transfer learning with a model pre-trained on the ImageNet dataset [3] to accelerate training and improve the model performance. The weights from the pre-trained model are applied to our model, but are not frozen and are subject to change during the training process. This is due to ImageNet having little connection to facial image classification.

The training procedure is outlined in the following steps. On the top layer, we implement global average pooling to average each feature map and feed it directly to the softmax layer [4]. In addition, a dropout is also added as a regularizer to prevent overfitting. The model is compiled using stochastic gradient descent with its learning rate optimized by the Adam algorithm [5] and sparse categorical cross entropy to calculate the loss. We train the model using up to 1000 epochs, but an early stop callback discontinues the training when validation accuracy doesn't reach its recorded maximum. In addition, a model checkpoint callback is used to save the best training progress.

## Weighted

Our first ablation study attempts to correct the imbalanced dataset by implementing class weights. We calculate the weights for each class by dividing the total number of pairs by the number of pairs of that class, then divide that by the number of classes. The weights are applied to all of the classes, adding importance to the under-represented classes. The model is then trained using these new class weights, following the same procedure as the baseline model.

## Oversampling

Another approach to address the imbalanced dataset was to oversample the minority classes. The oversampling process involves randomly sampling image pairs for each class to make the number of total pairs per class equal to the number of mother-son pairs at 46,969, so all of the other classes are then sampled to have that same number of pairs. This new dataset of classes of equal size is then used as the training data and the model is trained following the same procedure as the baseline model.

## Augmentation

We attempt another ablation study through the augmentation of images within the oversampled training dataset. The data is randomly selected to be augmented and the changes include flipping images horizontally and slightly rotating them by $0.06(2\pi)$ radians. Applying this method serves the purpose of providing more data for the model to train on. We then train the augmented model using the same process as the baseline model.

## Siamese Network

Lastly, we try an ablation study using a Siamese network. This architecture contains two embedded models with the same configuration for parameters and weights [6]. The oversampled dataset is preprocessed similarly to the previous models, but instead of taking two-channel tensors on one model, it takes one-channel tensors on each of the embedded models. The outputs from both of these embedded models are concatenated together, which then is further trained using a similar procedure to the baseline model.

# 4    Results

The Results section is structured by displaying plots of the accuracy and loss for both training and validation data as sub-figure (a) in the pair of figures. The accuracy is in decimal format, not percent. In sub-figure (b) the confusion matrix and corresponding heat map for the validation fitting are visualized. The confusion matrix displays the accuracy for each class of pairs, where each row represents instances of the actual classes, from 0 to 11 corresponding to the

types of pairs possible [Figure 4]. Each column represents the instances of the predicted class, also from 0 to 11 with the same correspondence [Figure 4]. The diagonals, starting from top left to bottom right, represent the number of times a class was properly fit, as the column number equals the row number. When the row number does not equal the column number, that shows the number of times that incorrect fitting was made. The heat map below the matrix visualizes this confusion matrix, with the colors representing the percent over the total amount of pairs for that actual class.

```
'fs': 0,
'ms': 1,
'bb': 2,
'fd': 3,
'md': 4,
'ss': 5,
'sibs': 6,
'gfgs': 7,
'gmgs': 8,
'gmgd': 9,
'gfgd': 10]
```

Figure 4: The classification names for each type of pair, along with the number that corresponds to the type of pair for labeling.

## Baseline

The results of the Baseline model are shown in Figure 5. The trained model is saved as "baseline.h5" in the "/Model" directory of the Google Drive.

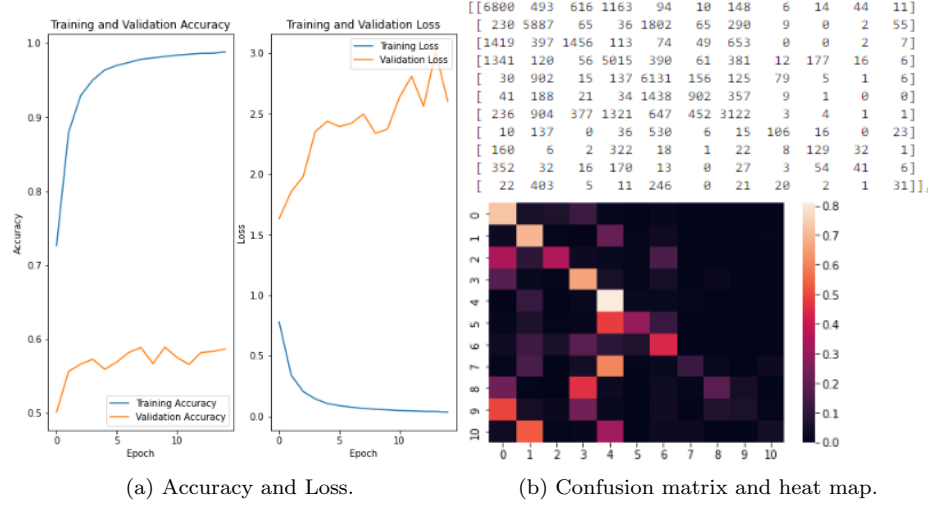(a) Accuracy and Loss.

(b) Confusion matrix and heat map.

Figure 5: Results for the Baseline model.

For the Baseline model, the training accuracy steadily improved as Epochs progressed, starting at 73% and maximizing at about 99%, with loss decreasing from approximately 0.78 to about 0.03. Validation accuracy increased from 50% to around 59%, with loss also increasing nearly linearly as the epochs progressed, from 1.6 to 2.6. The confusion matrix and heat map show a trend, with the lower numbered classes (0-6) having better fit results than the higher number classes (7-11). The lower number classes correspond to first-degree relationship pairs, the parents and siblings, while the higher number classes represent the second-degree pairs, the grandparent to grandchild pairs.

## Weighted

The results of the Weighted model are shown in Figure 6. The trained model is saved as "weighted.h5" in the "/Model" directory of the Google Drive.

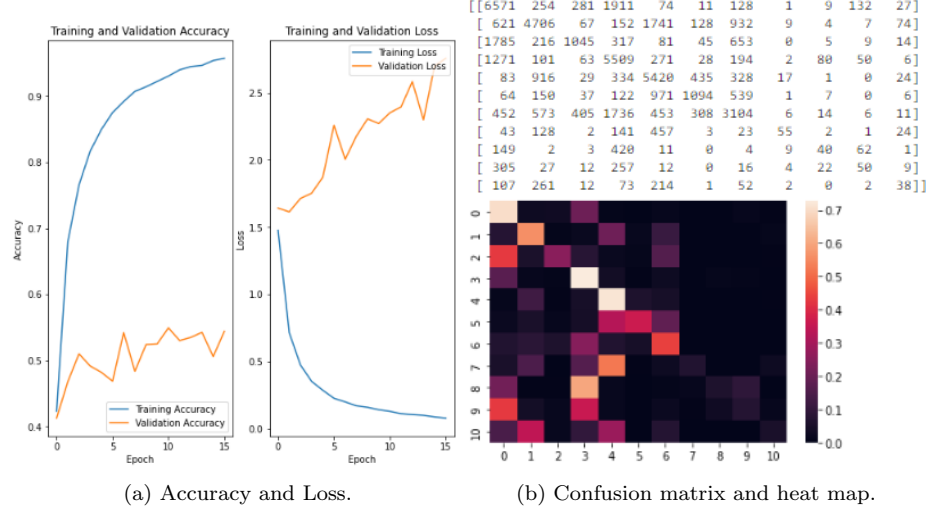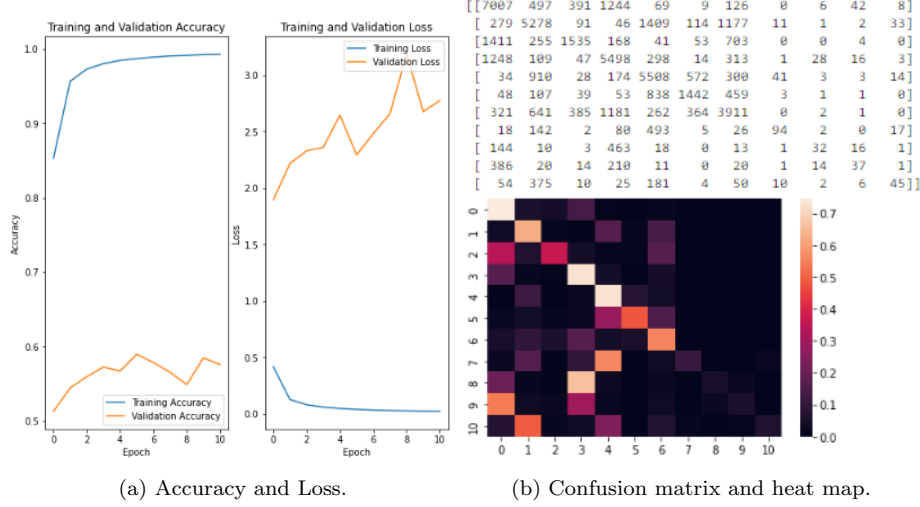(a) Accuracy and Loss.

(b) Confusion matrix and heat map.

Figure 6: Results for the Weighted model.

For the Weighted model, the training accuracy improved from about 42% to approximately 96%, with loss decreasing from 1.4 to 0.08. The validation accuracy peaked at 54%. The validation loss increased to a maximum of 2.8. The confusion matrix and heat map show a similar trend to that of the baseline, with the first-degree relationships having significantly better accuracy than the second-degree relationships.

## Oversampling

The results of the Oversampling model are shown in Figure 7. The trained model is saved as "oversample.h5" in the "/Model" directory of the Google Drive.

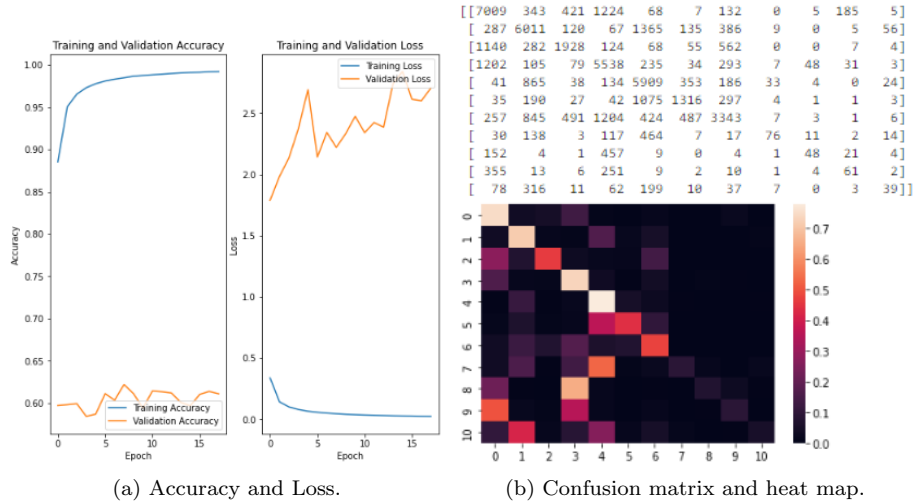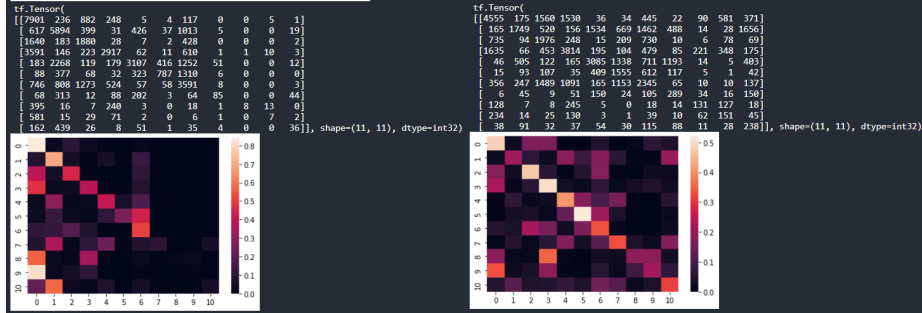(a) Accuracy and Loss.

(b) Confusion matrix and heat map.

Figure 7: Results for the Oversampling model.

For the Oversampling model, the training accuracy improved from approximately 85% to about 99%. The training loss decreased from about 0.42 to 0.02. The validation accuracy increased from approximately 51% to 58.5%, while its loss increased from approximately 1.9 to 2.8. The confusion matrix and heat map show a similar trend to that of the weighted model, with the first-degree relationships having better accuracy than the second-degree relationships.

## Augmentation

The results of the Augmentation model are shown in Figure 8. The trained model is saved as "augment_For_Sponsor.h5" in the "/Model" directory of the Google Drive.

(a) Accuracy and Loss.  (b) Confusion matrix and heat map.

Figure 8: Results for the Augmentation model.

For the Augmentation model, the training accuracy increased from slightly below 80% to around 91%, while its loss decreased from approximately 0.5 to below 0.25. The validation accuracy peaked at about 61%, while the loss increased as the Epochs progressed, from about 1.6 to 2.1. The confusion matrix and heat map show the same trend as the previous models, with the first-degree relationships having better accuracy than the second-degree relationships.

## Siamese Network

The results of the Siamese Network are shown in Figure 9. The trained model is saved as "siamese.h5" in the "/Model" directory of the Google Drive.
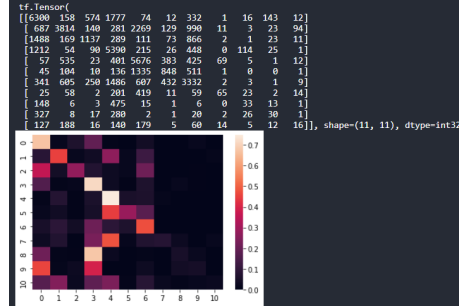


Figure 9: Results of incomplete Siamese model training.

For the Siamese Network, the training process was taking significantly longer for each epoch, while the performance was not better than our previous model, so we discontinued after the fifth epoch.

## One Epoch test results

(a) Baseline confusion matrix and heat map.  (b) Weighted confusion matrix and heat map.



(c) Oversampled confusion matrix and heat map.

Figure 10: Results from running only one epoch.

# 5    Conclusion

In this project, we were able to leverage, train, and finely tune a CNN model on the Families in the Wild dataset to determine the familial relationships between sets of two people. The imbalanced nature of the dataset prompted us to experiment with adjustments to the baseline model, using weighting and oversampling. In an attempt to improve the accuracy, we also implement augmentation to the oversampled baseline. In addition, we attempted a Siamese Network to compare its results with the previous models. The results of these models differ slightly, with Augmentation proving to be the most accurate with validation accuracy reaching 62.2%, followed by Oversampling at 59.0% accuracy, Baseline at 58.9%, and Weighted at 55.0%. We found that the Siamese model takes significantly longer to train on each epoch, due to the two models are being simultaneously trained, while producing similar results to the rest of the models. Out of the different methods, Augmentation returns the highest accuracy because it provides more randomness to the oversampled dataset.

For all of the models, we observed that as validation accuracy increased, so did its loss. We hypothesize that this is because the four second-generation

11

classes aren't as recognizable to the model, due to the imbalanced data providing fewer instances of these relationships to train on. Even when using oversampling and weighting, the model's attempts to minimize the loss in training lead to overfitting on these last 4 classes, while still learning new things about the other classes. We developed this hypothesis based on tests done using only one epoch on each of these different models [Figure 10]. In the first epoch, the weighted results [Figure 10b] fit much better for the minority classes than either the baseline [Figure 10a] or oversampled [Figure 10c] models. However, the oversampled model does perform better than the baseline on the minority classes, while achieving the same accuracy. This reinforces the notion that as more training is conducted, the fit on the minority classes becomes worse.

A potential future approach to achieve better accuracy would be to use a Generative Adversarial Network [7]. This would be used instead of oversampling, as it generates more data for the minority groups. Another improvement we could envision would be to include metadata as a class, such as when the picture was taken or the age of the individual in the portrait, as this could assist the model with instances in the dataset of pictures of the same person from different decades.

# References

[1]  Joseph P. Robinson et al. *Families In the Wild A Large-Scale Kinship Recognition Image Database*. URL: https://web.northeastern.edu/smilelab/fiw/index.html. (accessed: 11/3/2021).

[2]  Mingxing Tan and Quoc V. Le. "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks". In: *International Conference on Machine Learning* (2019). DOI: https://arxiv.org/abs/1905.11946.

[3]  J. Deng et al. "ImageNet: A Large-Scale Hierarchical Image Database". In: (2009).

[4]  Min Lin, Qiang Chen, and Shuicheng Yan. "Network In Network". In: (2014), p. 4. arXic: 1312.4400 (cs.NE).

[5]  Diederik P. Kingma and Jimmy Ba. "Adam: A Method for Stochastic Optimization". In: (2017). arXiv: 1412.6980 [cs.LG].

[6]  Davide Chicco. "Siamese Neural Networks: An Overview". In: (2021). Ed. by Hugh Cartwright, pp. 73–94. DOI: 10.1007/978-1-0716-0826-5_3. URL: https://doi.org/10.1007/978-1-0716-0826-5_3.

[7]  Antonia Creswell et al. "Generative Adversarial Networks: An Overview". In: *IEEE Signal Processing Magazine* 35.1 (Jan. 2018), pp. 53–65. ISSN: 1053-5888. DOI: 10.1109/msp.2017.2765202. URL: http://dx.doi.org/10.1109/MSP.2017.2765202.