

# **Automatic Face Understanding: Recognizing Families in Photos**

A Dissertation Presented

by

**Joseph P Robinson**

to

**The Department of Electrical and Computer Engineering**

in partial fulfillment of the requirements

for the degree of

**"Doctor of Philosophy"**

in

**Computer Engineering**

**Northeastern University  
Boston, Massachusetts**

February 11, 2023

*To my mom, who has always shown me the world as a world of opportunities.*

# Contents

<b>List of Figures</b>	<b>vi</b>
<b>List of Tables</b>	<b>xiii</b>
<b>List of Acronyms</b>	<b>xv</b>
<b>Acknowledgments</b>	<b>xvii</b>
<b>Abstract of the Dissertation</b>	<b>xxi</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Objective and Significance . . . . .	2
1.1.1 Scope . . . . .	2
1.2 Contributions . . . . .	3
1.2.1 Organization . . . . .	3
1.2.2 Publications . . . . .	4
1.2.3 Service . . . . .	6
<b>I Preliminaries</b>	<b>8</b>
<b>2 Automatic Facial Recognition</b>	<b>9</b>
2.1 Overview . . . . .	9
2.2 Traditional Methods . . . . .	12
2.2.1 Eigenfaces . . . . .	12
2.2.2 Fisherfaces . . . . .	17
2.2.3 Local Binary Pattern Histograms (LBPH) . . . . .	18
2.2.4 Results and analysis . . . . .	19
2.3 Modern-day, data driven deep learning . . . . .	20
2.3.1 Loss functions . . . . .	22
<b>3 Face Detection</b>	<b>24</b>
3.1 Overview . . . . .	24
3.2 Research Contributions . . . . .	25

3.3	Background Information . . . . .	27
3.3.1	Landmark localization . . . . .	27
3.3.2	GANs . . . . .	28
3.4	Laplace Landmark Localizer . . . . .	29
3.4.1	Fully Supervised Branch . . . . .	29
3.4.2	Unsupervised Branch . . . . .	31
3.4.3	Training . . . . .	32
3.4.4	Implementation . . . . .	33
3.5	Experiments . . . . .	34
3.5.1	Metric . . . . .	35
3.5.2	300W + MegaFace . . . . .	35
3.5.3	The Annotated Facial Landmarks in the Wild (AFLW) dataset . . . . .	39
3.5.4	Ablation Study . . . . .	39
3.6	Discussion . . . . .	40
<b>II</b>	<b>Visual Kinship Recognition of Families In the Wild</b>	<b>41</b>
<b>4</b>	<b>Visual Kinship Recognition</b>	<b>42</b>
4.1	Overview . . . . .	42
4.2	Background Information . . . . .	47
4.2.1	The evolution of the problem . . . . .	48
4.2.2	Humans recognizing kinship in photos . . . . .	49
4.3	Visual Kinship Problems . . . . .	51
4.3.1	Kinship verification . . . . .	52
4.3.2	Family classification . . . . .	53
4.3.3	Tri-subject verification . . . . .	54
4.3.4	Search and retrieval . . . . .	55
4.3.5	Multi-modal data . . . . .	55
4.3.6	Kin-based facial synthesis . . . . .	56
<b>5</b>	<b>Families In the Wild (FIW)</b>	<b>57</b>
5.1	Overview . . . . .	57
5.2	Related Works . . . . .	59
5.2.1	Related Databases . . . . .	59
5.2.2	Automatic Kinship Recognition . . . . .	59
5.2.3	Deep Kinship Recognition . . . . .	60
5.2.4	Semi-Automatic Image Tagging & Data Exploration . . . . .	60
5.3	Families in the Wild (FIW) Database . . . . .	61
5.3.1	FIW v0.1 . . . . .	61
5.3.2	Data Preparation . . . . .	63
5.3.3	Extending FIW . . . . .	64
5.4	Semi-Supervised Face Clustering . . . . .	73
5.4.1	Objective Function . . . . .	74
5.4.2	Solution . . . . .	74

5.5	Experiments . . . . .	75
5.5.1	Experimental Setting . . . . .	76
5.5.2	Kinship Verification . . . . .	78
5.5.3	Family Classification . . . . .	80
5.5.4	Proposed Semi-Supervised Clustering . . . . .	81
5.5.5	Transfer-Learning Experiment . . . . .	82
5.5.6	Human assessment using FIW . . . . .	82
5.5.7	Discussion . . . . .	85
5.6	Data challenges and incentives . . . . .	86
5.7	Experimental . . . . .	87
5.7.1	Task Evaluations, Protocols, Benchmarks . . . . .	87
5.8	Methodologies . . . . .	91
5.8.1	Traditional approaches . . . . .	93
5.8.2	Deep learning approaches . . . . .	95
5.8.3	Generative modeling approaches . . . . .	105
<b>6</b>	<b>FIW-MM</b>	<b>106</b>
6.1	Overview . . . . .	106
6.2	Related Work . . . . .	109
6.2.1	Kinship recognition . . . . .	109
6.2.2	Audio-visual data . . . . .	111
6.3	The FIW-MM Database . . . . .	112
6.3.1	Specifications . . . . .	113
6.3.2	Data pipeline . . . . .	113
6.4	Problem Definitions and Protocols . . . . .	117
6.4.1	Kinship verification . . . . .	119
6.4.2	Search & retrieval (missing child) . . . . .	121
6.5	Benchmarks . . . . .	122
6.5.1	Methodology . . . . .	122
6.5.2	Results . . . . .	125
6.5.3	Discussion . . . . .	125
6.6	Future Work . . . . .	126
6.7	Conclusion . . . . .	128
<b>III</b>	<b>Post Processing</b>	<b>129</b>
<b>7</b>	<b>Kinship Recognition - State of Technology</b>	<b>130</b>
7.1	Overview . . . . .	130
7.2	Technical Challenges . . . . .	130
7.2.1	Current limitations of SOTA . . . . .	131
7.2.2	The nature . . . . .	132
7.2.3	The environment . . . . .	132
7.2.4	The data and its distribution . . . . .	133
7.3	Applications . . . . .	134

7.4	Discussion . . . . .	136
<b>8</b>	<b>Bias in Face Recognition</b>	<b>137</b>
8.1	Overview . . . . .	137
8.1.1	Organization . . . . .	142
8.2	Related Work . . . . .	142
8.2.1	Bias in machine learning . . . . .	142
8.2.2	Bias in facial recognition . . . . .	143
8.2.3	Human bias in machine learning . . . . .	145
8.2.4	Imbalanced data and data problems in FR . . . . .	146
8.2.5	Domain adaptation and feature alignment . . . . .	147
8.2.6	Protecting demographic information in face recognition (FR) . . . . .	147
8.3	Balanced Faces In the Wild (BFW) . . . . .	148
8.3.1	The data . . . . .	149
8.3.2	Problem formulation . . . . .	151
8.3.3	Human assessment . . . . .	151
8.4	Methodology . . . . .	152
8.4.1	Problem statement . . . . .	153
8.4.2	Proposed framework . . . . .	156
8.5	Results and Analysis . . . . .	158
8.6	Experiments . . . . .	162
8.6.1	Common settings . . . . .	163
8.6.2	Debias experiment . . . . .	163
8.6.3	Privacy preserving experiment . . . . .	165
8.6.4	Ablation study . . . . .	169
8.7	Discussion . . . . .	169
<b>9</b>	<b>Discussion</b>	<b>171</b>
9.1	Broader impacts . . . . .	171
9.2	Future work . . . . .	172
9.3	Conclusion . . . . .	174
<b>Bibliography</b>		<b>175</b>

# List of Figures

2.1	<b>A generic verification system.</b> Regardless of the network specifications ( <i>i.e.</i> , independent of layer counts, layer sizes, type of metric set topmost). The aim is to map two image inputs to a single logistic value. Input images ( <i>i.e.</i> , $I^i$ and $I^j$ ) must then assume the same size as that of the input layer of the model. Then, the learned representation of the $i$ -th and $j$ -th faces are output from the last layer of the network prior to the classification layer. An aspect of the model shown is the mechanism used to fuse ( <i>i.e.</i> , single logit value output from $\mathcal{C}$ , which was fed the features of the $i$ -th and $j$ -th face encoding). The framework is inherently boolean, as the task is to map a sample pair ( <i>i.e.</i> , pair of faces) to a boolean class label ( <i>i.e.</i> , 1 if <i>genuine</i> and 0 for <i>impostor</i> ). . . . . .	10
2.2	<b>Generic recognition system.</b> Faces are often stored as encodings in a database through <i>enrollment</i> . <i>Recognition</i> is then to compare the encoding of an input face to those in the database. . . . .	11
2.3	<b>Illustration of the Face Space, which spans the area of the light-blue rectangle.</b> Note that the blue spheres in the Face Space ( $\Omega_1$ , $\Omega_2$ , $\Omega_3$ ) describe a face of a particular person, <i>i.e.</i> , used to identify or verify individual instances. Beyond that, but within the bounds of the Face Space, are varying faces of unknown type, <i>i.e.</i> , could be used, for instance, to detect any face. . . . .	12
2.4	<b>By taking the difference between each face and the mean of all faces <i>Eigenfaces</i> are normalized.</b> Results in images of facial structure, and resembling a ghost– some call these Ghostfaces, it is most common, and agreeably more appropriate, to refer to them as <i>Eigenfaces</i> . . . . .	13
2.5	<b>Eigenfaces for the same face under different lighting conditions.</b> The original face images and their <i>Eigenface</i> equivalent are shown on top row and bottom row, respectively. Notice the variation between <i>Eigenfaces</i> . . . . .	13
2.6	<b>Illustration depicting the process of obtaining LBH features from an image for a single pixel.</b> With every pixel processed, a histogram is generated to represent each image. . . . .	14
2.7	<b>LBPH of the same face under different lighting conditions.</b> Notice the light invariance that is inherited with this feature. . . . .	15
2.8	<b>Accuracy measure.</b> Correct (%) as function of training samples count per class for the AT&T dataset (a). Correct (%) as function of training samples count per class for the Yale B Extended dataset (b). . . . .	16

2.9	<b>Performance measure.</b> Training time as function of training samples count per class for the AT&T dataset (Left). Training time as function of training samples count per class for the Yale B dataset (Right).	17
2.10	<b>Learned spaces (visualization from [1]).</b> These schematics are derived from the respective loss function— <i>left-to-right</i> : traditional softmax, NSL, ArcFace, and LMCL.	22
3.1	<b>Problem statement.</b> Heatmaps generated by SAM-based models (middle block) and the proposed LaplaceKL (right block), each with heatmaps on the input images (left) and a zoomed-in view of an eye region (right). These heatmaps are confidence scores ( <i>i.e.</i> , probabilities) that a pixel is a landmark. Softargmax-based methods generate highly scattered mappings (low certainty), while the same network trained with our loss is concentrated ( <i>i.e.</i> , high certainty). We further validate the importance of minimizing scatter experimentally (Table 3.2). Best if viewed electronically.	25
3.2	<b>Our semi-supervised framework for landmark detection.</b> The labeled and unlabeled branches are marked with <b>blue</b> and <b>red</b> arrows, respectfully. Given an input image, generator ( $G$ ) produces $K$ heatmaps, one for each landmark. Labels are used to generate real heatmaps as $\omega(\mathbf{s}^l)$ . $G$ produces fake samples from unlabeled data. Source images are concatenated on heatmaps and passed to discriminator ( $D$ ).	29
3.3	<b>Random samples (300W).</b> Heatmaps predicted by our LaplaceKL+D(70K) (middle, <i>i.e.</i> , L-KL+D(70K)) and softargmax+D(70K) (right, <i>i.e.</i> , SAM+D(70K)) alongside face images with ground-truth sketched on the face (left). For this, colors were set by value for the $K$ heatmaps generated for each landmark ( <i>i.e.</i> , range of $[0, 1]$ as shown in color bar), and then were superimposed on the original face. Note that the KL-divergence loss yields predictions of much greater confidence and, hence, produced separated landmarks when visualized heatmap space. In other words, the proposed has minimal spread about the mean, as opposed to the softargmax-based model with heatmaps with individual landmarks smudged together. Best viewed electronically.	36
3.4	<b>Qualitative results.</b> Random samples of landmarks predicted using LaplaceKL (white), with the ground truth drawn as line segments ( <b>red</b> ). Notice the predicted points tend to overlap with the ground-truth. Best viewed in color. Zoom-in for greater detail.	37
3.5	<b>Ablation.</b> Results of ablation study on LaplaceKL.	39

4.1	<b>A decade of research in visual kinship recognition.</b> The timeline shows correlations between the data resources ( <i>below timeline</i> ) and citation metrics and events indicating the amount of research impact ( <i>above timeline</i> ). We built a pipeline to scrape the data needed for the plots above: (1) <i>Publish or Perish</i> [2] was installed on a Mac Book Pro to gather metadata for publications from various sources ( <i>i.e.</i> , Google Scholar, Cross Ref, and Scopus) into a CSV file; (2) metadata in CSV was parsed into a BIB file using Python; (3) <i>Mendeley Reference Manager</i> was used to automatically detect duplicates while keeping as much information as possible by merging reference listings; (4) queried Google Scholar for all <i>Related Works</i> and <i>Cited By</i> using PyPi’s scholarly ( <a href="https://pypi.org/project/scholarly/">https://pypi.org/project/scholarly/</a> ), which extended the paper-pile and increased the amount of metadata available from the richer metadata accessible using scholarly ( <i>e.g.</i> , paper abstracts); (5) we clustered the documents by abstract via term frequency-inverse document frequency learning (TF-IDF) [3]. The clusters were high in recall, as true clusters were a majority of papers on kinship recognition in multimedia: this reduced the burden of manual inspection of hundreds of thousands to thousands. It is important to note that only citation metrics were considered, leaving out other factors of impact like the <i>number of times tweeted</i> , <i>Github stars</i> , and other indicators of impacting research. . . . .	44
4.2	<b>Workflow to scrape publication metadata for Figure 4.1.</b> From <i>Publish or Perish</i> [2], we queried Scholar for <i>Related works</i> and <i>Cited by</i> , increasing the size of our list nearly 20-fold. Mendeley merged duplicates, while keeping as much information as possible. Applied Natural Language Processing (NLP) to cluster relevant documents. . . . .	46
4.3	<b>Samples used for human evaluation.</b> Each column displays pairs most commonly marked correctly and incorrectly, and in cases where the correct answers were true and false. Specifically, true-positive (TP), false-positive (FP), true-negative (TN), and false-negative (FN) are displayed, respectively. Each of these pairs was properly classified by the fine-tuned CNN. . . . .	51
4.4	<b>Visual kin-based discriminate tasks for the <i>Families In the Wild</i> (FIW) dataset.</b> Robinson <i>et al.</i> posed problems of verification ( <i>i.e.</i> , one-to-one) [4] and family classification ( <i>i.e.</i> , one-to-many) [5, 6], along with more recently supporting tri-subject verification ( <i>i.e.</i> , one-to-two) and search & retrieval for “missing” children ( <i>i.e.</i> , many-to-many) [7]: the FIW database supported the aforementioned tasks, while the annual data challenge <i>Recognizing Families In the Wild</i> (RFIW) was, and continues to be, motivated by promotional purposes. The most recent data challenge supported include three of the four shown here, as family classification was found to carry less potential for practical use-cases, while the others were done using three data splits disjoint in terms of family labels (Table 5.9, 5.11, 5.13). Protocols and benchmarks for each view are described in [7]. Best viewed electronically. . . . .	52
4.5	<b>Generic Siamese network.</b> Approaches tend to follow the Siamese model, differing in method of fusion. Specifically (from <i>top-to-bottom</i> ), an image pair shot $x_1$ and $x_2$ . . . . .	54
5.1	<b>Sample family photos from FIW.</b> Randomly picked 8 / 1,000 families. . . . .	58
5.2	<b>Samples of eleven pair types of FIW.</b> Each type is of a unique pair randomly selected from a set of diverse families, while four faces of each individual depict age variations. . . . .	61

5.3	<b>Database statistics.</b> Horizontal and vertical axes represent counts for photos and faces per family, respectively. Bubble size and color represent counts for members and average faces per member, respectively. . . . .	63
5.4	<b>Visual of the label types of FIW, Family-level (FID) and Photo-level (PID).</b> FID has individual family member (MID) and relationship information. PIDs contain information of MIDs + their locations in photos. . . . .	64
5.5	<b>Semi-automatic labeling pipeline.</b> <i>Data Collection.</i> Photos and text metadata were collected for underrepresented families in FIW and assigned unique IDs ( <i>i.e.</i> , PIDs). Each new member requires at least 1 profile picture ( <i>e.g.</i> , Brandon in $PID_1$ ) to add to known labels. <i>Data Preparation.</i> With the existing FIW labels, we next aim to increase the amount, both in labeled faces and member labels, using multiple modalities—names in metadata and scores of Support Vector Machines (SVMs) were used to automatically label some unlabeled data—face-name pairs were assumed labeled for cases of high confidence. Starting from profile pictures ( <i>i.e.</i> , 1 face, 1 name) and working towards less trivial scenarios ( <i>e.g.</i> , 3 faces and 2 names, with 2 faces from 1 member at different ages, like in $PID_3$ ). This step adds to the amount of side information used for clustering. <i>Label Generation.</i> Label proposals for remaining unlabeled faces were generated using the proposed semi-supervised clustering model that leverages labeled data as side information to better guide the process. <i>Label Validation.</i> A GUI designed to validate clusters and ensure clusters matched the proper labels. . . . .	65
5.6	<b>Visualization depicting family structure and photos of the Royal Family.</b> There are several members in the tree (top) and many photos in total (bottom). . . . .	67
5.7	<b>Family photo montage.</b> Samples photos for 8 of 1,000 families in FIW. . . . .	68
5.8	<b>Bruce Lee family tree before (top) and after (bottom) extension.</b> The photos in the middle were added to existing photos using our proposed semi-automatic labeling scheme. This increased both the samples per members and the total number ( <i>i.e.</i> , from just 3 to 10 members). . . . .	69
5.9	<b>Sample family of FIW [6].</b> Faces and relationships of the American Football family, the Gronkowski's ( <i>Top</i> ). The montage shows less than half of all photos for respective family. Photo types are various, spanning profile faces ( <i>top</i> ) to images of different subgroups of family members. Furthermore, samples capture different times of life. Note, crops were made to fit montage ( <i>Bottom</i> ). . . . .	72
5.10	<b>Relationship type specific ROC curves.</b> Notice the fine-tuned Sphereface dominates, while the sample counts for the <i>grandparent-grandchild</i> were less as indicative of the jagged curves. . . . .	77
5.11	<b>Results for clustering families using different amounts of side information.</b> As clearly depicted, our method obtains the top performance. Moreover, a distinct increase in NMI for our method is shown with an increase in the amounts of side information. . . . .	83
5.12	<b>Box plot for humans on kinship verification.</b> <i>Case 1:</i> Relationship type dependent evaluations. <i>Case 2:</i> Evaluations with type unspecified. . . . .	84

5.13	<b>Kinship verification (T1) sample pairs.</b> Sample pairs with similarity scores near the threshold ( <i>i.e.</i> , hard (H) samples), along with highly confident predictions ( <i>i.e.</i> , easy (E) samples) in verification task. . . . .	92
5.14	<b>Qualitative analysis of T1.</b> Samples of each relationship type that all of the teams either got correct (100%) or mostly not (20%) for the eleven pair types of FIW and NON-KIN. . . . .	93
5.15	<b>Triplets with extreme scores (<i>i.e.</i>, correct and incorrect).</b> Each show FMS (top rows) and FMD (bottom) for tri-subject (T2). . . . .	94
5.16	<b>Sample of T2.</b> Samples that all teams got correct (left) and mostly incorrect (right) for FMS (top rows) and FMD (bottom). . . . .	95
5.17	<b>Plot of face counts per family in test set of T3.</b> The probes have about 8 faces on average, while the number of family members in the gallery nears 20 on average, with a total average of 170 faces. . . . .	97
5.18	<b>T3 sample results (Rank 10).</b> Each query (row) has one or more faces, for the probe returns and ranks all samples in the gallery - here we show top 10. FP are labeled by <b>x</b> , while true matches list the relationship type in green: <b>P</b> for parent; <b>C</b> for child; <b>S</b> for sibling. . . . .	99
5.19	<b>Activations from mapping image-to-latent space (from [8]).</b> The salience mapped from the activation response and superimposed on the average face. Family101 dataset was used for this experiment [9]. The end result depicted here were dubbed the <i>genetic features</i> from latent space of a trained Gated autoencoder (AE). . . . .	100
5.20	<b>Model to synthesize children faces from a parent-pair (visualizations from [10]).</b> Notice that the output of encoder $E$ is the concatenation of features from prospective parents, the father $h_f$ and mother $h_m$ joined by $\oplus$ such that the two embeddings encoded by the Siamese network are fused ( <i>i.e.</i> , $2 * \mathbb{R}^d \rightarrow \mathbb{R}^{2d}$ ) before passed as input to the conditional adversarial AE (CAEE) model. . . . .	101
5.21	<b>Synthesized results from [10].</b> Columns correspond to families, with fathers on first row, mothers on second, and real and generated children on third row and bottom, respectively (a). See subcaption for specifics on age (b) and gender (c). . . . .	102
5.22	<b>Salience map per key-points [10].</b> Best viewed in color. . . . .	103
6.1	<b>Sample family of FIW in Multimedia (FIW-MM).</b> Top-to-bottom: <i>family-tree labels</i> show faces of members in the immediate family, with subjects of the same generations in the same row; <i>videos, audio, and contextual</i> exemplify sample video pairs of Dr. King Jr. and his daughter Andrea with tracklets of faces in the visual domain and audio data aligned frame-by-frame; <i>family photos</i> that contain Dr. Luther King Jr. randomly selected (note, cropped to fit); <i>faces</i> of Dr. King Jr. from adolescence-to-adulthood. Multiple faces are available for most subjects. Best viewed electronically. . . . .	108
6.2	<b>Automated labeling framework.</b> For each of the 1,000 families, there are a set of $K$ members. For this, the template of a member consists of all media available. Tag numbers 1-6 correspond to sections in Section 6.3.2 . . . . .	111

6.3	<b>Plotted results.</b> Shown are score fusion (SF) and feature fusion (FF) as late and early fusion methods, respectively. Included are still-images $S$ , video clips $V$ , and audio segments $A$ , with still-images and video were fused $S + V$ , and also still-images, audio, and video were fused $S + V + A$ . Clearly, both tasks benefit from early fusion: the DET curve (left) summarizes the verification task by plotting FNR as a function of FPR ( <i>i.e.</i> , lower is better); search and retrieval is summarized as a CMC (right) by showing the accuracy as a function of rank ( <i>i.e.</i> , higher is better). . . . .	121
6.4	<b>Random hard sample.</b> Template of a true FS pair that was incorrectly classified using score fusion, but correct for TA ( <i>i.e.</i> , feature fusion). Here only a single face is available for the father (left), while all instances of the son are at a young age (right). . . . .	125
6.5	<b>Hard sample pair.</b> A true MS pair incorrectly classified using late fusion, while correctly identified as KIN when using the early fusion via TA. Challenges here are in the young age of the son and the majority of the faces of the mother occluded by sunglasses ( <i>i.e.</i> , score fusion puts equal weight on all samples, where TA learns to better discriminate). . . . .	127
7.1	<b>Sample faces synthesized to improve predictive power for faces of elderly adults (visualization from [11]).</b> Two models ( <i>i.e.</i> , one per gender) were trained to synthesize input faces as younger– male fathers ( <i>i.e.</i> , rows 1-2) and female mothers ( <i>i.e.</i> , rows 3-4), the top sample is the original and the generated is below. . . . .	132
7.2	<b>‘From same photograph’ (FSP) data (modified from [12]).</b> Data tagged via constraints, <i>must</i> and <i>cannot-link</i> : $\approx 1M$ data points scraped from the web via 125 non-kin queries ( <i>e.g.</i> , school student, sports team). . . . .	133
8.1	<b>Depiction of the biometrics.</b> The SDM shows the sensitivity related to a single threshold $t_g$ ( <i>top-left</i> ). The area to the right of the threshold considers all accepted pairs, both correctly and incorrectly predicted. True Acceptance Rate (TAR) as a function of False Acceptance Rate (FAR) is a common way to report ratings for given false-rates ( <i>top-right</i> ). Equally common in FR is the trade-off between false-negative rate (FNR) and FAR ( <i>bottom-left</i> and <i>bottom-right</i> ). . . . .	138
8.2	<b>BFW dataset.</b> Average face per subgroup: <i>top-left</i> : the entire Balanced Faces In the Wild (BFW); <i>top-row</i> per ethnicity; <i>left-column</i> : per gender. The others represent the ethnicity and gender, respectively. Table 8.1 defines the acronyms of subgroups. . . . .	144
8.3	<b>Signal detection model (SDM) across subgroups.</b> Scores of <i>imposters</i> have medians $\approx 0.3$ but with variations in upper percentiles; <i>genuine</i> pairs vary in mean and spread ( <i>e.g.</i> , AF has more area of overlap). A threshold varying across different subgroups yields a constant FAR. . . . .	145
8.4	<b>Detection error trade-off (DET) curves.</b> <i>Top-left</i> : per gender. <i>Top-right</i> : per ethnicity. <i>Bottom</i> : per subgroup ( <i>i.e.</i> , combined). Dashed line shows about $2\times$ difference in false-positive rate (FPR) for the same threshold $\theta_{const}$ . FNR is the match error count (closer to the bottom is better). . . . .	150
8.5	<b>Confusion matrix.</b> Error (Rank 1, %) for all BFW faces versus all others. Errors concentrate intra-subgroup - consistent with the SDM (Figure 8.3). Although subgroups are challenging to define, this shows the ones chosen are meaningful for FR. . . . .	152

8.6	<b>BFW statistics (i.e., pixel counts).</b> Histogram of image areas in pixels (blue plot). The orange curve shows the cumulative count of images up to a given area. . . . .	153
8.7	<b>Samples of BFW.</b> Per subgroup: the 25 samples for a random subject are shown. . . . .	154
8.8	<b>Debiasing framework.</b> The framework used to learn a projection that casts facial encodings to a space that (1) preserves identity information (i.e., $C_{ID}$ ) and (2) removes knowledge of subgroup (i.e., $C_{ATT}$ ). The benefits of this are two-fold: ability to verify pairs of faces fairly across attributes and an inability to classify attribute for privacy and safety purposes. Note, that the <i>gradient reversal</i> [13] flips the sign of the error back-propagated from $C_{ATT}$ to $M$ by scalar $\lambda$ during training. . . . .	155
8.9	<b>DET curves for different CNNs.</b> FNR (%) (vertical) vs FPR (horizontal, log-scale) for VGG2 [14] models with different backbones (VGG16 [15], Resnet50 [16], SENet50 [17], in that order). Lower is better. For each plot, <i>White-Male</i> (WM) is the top-performer, while <i>Asian-Female</i> (AF) is the worst. The ordering of the curves is roughly the same for each backbone. . . . .	158
8.10	<b>Human assessment (qualitative).</b> ✓ for <i>match</i> ; ✗ for <i>non-match</i> . Accuracy scores shown as bar plots. Humans are most successful at recognizing their own subgroup, with a few exceptions (e.g., bottom). . . . .	159
8.11	<b>Percent difference from intended FPR per subgroup.</b> <i>Top:</i> global threshold ( $t_g$ ) yields a FPR that spans up to 200% the intended (i.e., WM for 1e-4); the F subgroups tend to perform worse than intended for all, while M's overshoot the intended performances besides IM at FPR=1e-4. <i>Bottom:</i> Subgroup-specific (or optimal) thresholds $t_o$ reduces the difference closer to zero. Furthermore, the proposed method ( <i>middle</i> ), which does not assume knowledge of attribute information at inference like for $t_o$ , clearly mitigates the issue of the inconsistencies in the true versus reported FPR. Similar to the results in Table 8.5, the variations are nearly halved: the percent difference for subgroups is more balanced using the adapted features versus the baseline. . . . .	162
8.12	<b>Sample pairs of hard positives.</b> Pairs incorrectly classified by the baseline and correctly matched by the proposed. . . . .	166
8.13	<b>Subgroup confusion matrix.</b> Comparison of accuracy in classifying and misclassifying the subgroups. Notice the (b) performs significantly worse than (a) as intended. 168	
8.14	<b>Accuracy on LFW benchmark.</b> The proposed approaches the performance of the baseline before dropping off. . . . .	170

# List of Tables

3.1	$12_{G_{h \times w \times n}} x_{E_{ID} - D_{ID}} \rightarrow K$	33
3.2	<b>Quantitative results.</b> Normalized Mean Square Error (NMSE) on Annotated Facial Landmarks in the Wild (AFLW) and 300W normalized by the square root of BB area and interocular distance, respectfully.	34
3.3	<b>Ablation.</b> NMSE on 300W (full set) for networks trained with fewer channels in each convolutional layer by 1/16, 1/8, 1/4, 1/2, and unmodified in size ( <i>i.e.</i> , the original) listed from left-to-right. We measured performance with a 2.8GHz Intel Core i7 CPU.	38
4.1	<b>Publicly available datasets for kinship recognition.</b> Each listed by the original name per reference. Kin-based image (or video) stats, which include the label types that support a specific evaluation metric and the respective state-of-the-art (SOTA) score. URLs to the project page of each data resource are included. Abbreviations used for <i>Stats</i> are for the family count ( <b>F</b> ), face count ( <b>f</b> ), number of unique people ( <b>P</b> ), sample count ( <b>S</b> ), image count ( <b>I</b> ), video count ( <b>V</b> ), and multimedia ( <b>MM</b> ). . . . . .	45
4.2	<b>KinWild benchmarks.</b> Results for KinWild I and II. . . . .	55
5.1	<b>Pairwise counts of FIW.</b> Notice FIW is first to provide Grandparent-Grandchild pairs. Table 5.2 further characterizes that data, and Figure 5.2 shows samples from it. . . . .	66
5.2	<b>Database counts and attributes.</b> Comparison of FIW with related datasets. . . . .	71
5.3	<b>Ethnicity distribution for the 1,000 of FIW.</b> Mix families contain >2 ethnicity ( <i>e.g.</i> , Bruce ( <i>Asian</i> ) and Linda ( <i>Caucasian</i> ) Lee with 2 children. . . . .	71
5.4	<b>Speedup analysis.</b> Previous (white) versus new (shaded) labeling processes compared in terms of inputs (keyboard and mouse clicks) and time (hours:minutes:seconds). . . . .	73
5.5	<b>Averaged verification accuracy scores (%) for 5-fold experiment on FIW.</b> Note that there was no family overlap between folds. . . . .	76
5.6	<b>Family classification results.</b> Accuracy scores (%) using 564 families. . . . .	80
5.7	<b>Transfer learning experiment.</b> Accuracy (%) for KinWild I & II. CNN fine-tuned on FIW top scorer. Note that these results were up-to-date when journal ( <i>i.e.</i> , [6]) was released, but is no longer. See Table 4.2 for most up-to-date scores. . . . .	82
5.8	<b>Face pair counts for human evaluation on kinship verification.</b> SIBS represents all siblings of the same generation, PC are parent-child, and GPGC are grandparent-grandchild. . . . .	85

5.9	<b>Kinship verification (T1) counts.</b> Number of unique pairs ( <b>P</b> ), families ( <b>F</b> ), and face samples ( <b>S</b> ), with an increase in counts and types since [5]. . . . .	86
5.10	<b>Kinship verification (T1) results.</b> Averaged verification accuracy scores of RFIW. . . . .	88
5.11	<b>Tri-subject verification (T2) counts.</b> No. pairs ( <b>P</b> ), families ( <b>F</b> ), face samples ( <b>S</b> ). . . . .	92
5.12	<b>Verification scores.</b> Results for tri-subject ( <i>i.e.</i> , T2). . . . .	94
5.13	<b>Tri-subject (T2) counts.</b> Individuals <b>I</b> , families <b>F</b> , face samples <b>S</b> . . . . .	96
5.14	<b>T3 results.</b> Performance ratings for SOTA methods. . . . .	98
6.1	<b>Database statistics.</b> Types are split based on the span in generation of the relationship. . . . .	113
6.2	<b>Task-specific counts.</b> For individuals ( <b>I</b> ), families ( <b>F</b> ), still-face images ( <b>S</b> ), video-clips ( <b>V</b> ), audio snippets ( <b>A</b> ), audio snippets ( <b>VA</b> ) in the set of probes ( <b>P</b> ), gallery ( <b>G</b> ), and in total ( <b>T</b> ). . . . .	115
6.3	<b>TAR at specific FAR.</b> Scores are for template-based settings: still-images only (left column), +videos (middle), and +video+audio (right). Higher is better. . . . .	117
6.4	<b>Identification results, with True Acceptances (TAs) highlighted.</b> Accuracy scores for different ranks are listed ( <i>i.e.</i> , higher is better). Also, MAP scores are provided for each. . . . .	120
8.1	<b>Database stats and nomenclature.</b> <i>Header:</i> Subgroup definitions. <i>Top:</i> Statistics of BFW. <i>Bottom:</i> Number of pairs for each partition. Columns grouped by ethnicity and then further split by gender. . . . .	139
8.2	<b>Data statistics, notation, and scores for subgroups of our BFW data.</b> <i>Top:</i> Specifications of BFW and subgroup definitions. <i>Middle:</i> Number of pairs. <i>Bottom:</i> Accuracy fo a global threshold $t_g$ , the value of the optimal threshold $t_o$ , and accuracy using $t_o$ per subgroup. Columns grouped by race and then further split by gender. Notice the inconsistent ratings across subgroups. . . . .	153
8.3	<b>Human assessment (quantitative).</b> Subgroups listed per row ( <i>i.e.</i> , human) and column ( <i>i.e.</i> , image). Note, most do the best intra-subgroup ( <b>blue</b> ), and second-best ( <b>red</b> ) intra-subgroup but inter-gender. WF performs the best; WF pairs are most correctly matched. . . . .	160
8.4	<b>BFW features compared to related resources.</b> Note, the balance across identity (ID), gender (G), and ethnicity (E). Compared with Demographic Pairs (DemogPairs), BFW provides more samples per subject and subgroups per set. Also, BFW uses a single resource, VGG2. Racial Faces in-the-Wild: (RFW); on the other hand, supports a different task ( <i>i.e.</i> , subgroup classification). Furthermore, RFW and FairFace focus on race-distribution without the support of identity labels. . . . .	161
8.5	<b>True Acceptance Rate (FAR) for various False Acceptance Rate (FAR).</b> TAR scores for a global threshold (top), the proposed debiasing transformation (middle), optimal threshold (bottom). Higher is better. The standard deviation from the average is shown to demonstrate the standard error comparing the reported ( <i>i.e.</i> , average) to the subgroup-specific scores. The proposed recovers most of the loss from using a global threshold rather than a per-subgroup threshold. . . . .	165
8.6	<b>Subgroup classification results.</b> The baseline and proposed are on the left and right columns, respectively. Note that the columns on the right have lower scores as intended. . . . .	167

# List of Acronyms

- AE** Auto Encoder. A model consisting of an encoder and a decoder module, where the encoder projects input signal to a latent, hidden state and the decoder reconstructs the original input.
- AP** Average Precision. A measure of the number of true-positives per total number of positives.
- BB** Brother-brother. Relationship type (pairwise).
- BFW** Balanced Faces in the Wild. A facial recognition dataset for measuring bias across different demographics (*i.e.*, ethnicity and gender).
- CMC** Cumulative matching characteristic [curve].
- DL** Deep learning.
- DET** Detection error trade-off [curve].
- FAR** False-acceptance rate.
- FD** Father-daughter. Relationship type (pairwise).
- FID** Family ID. Unique identifier assigned to each family of FIW.
- FIW** Families In the Wild. A large-scale dataset for recognizing family members in photos.
- FIW-MM** Families In the Wild in Multimedia. A large-scale dataset for recognizing family members in multimedia (*i.e.*, photos, video, audio, and text transcripts).
- FMD** father/mother-daughter. Relationship type (triplet).
- FMS** father/mother-son. Relationship type (triplet).
- FR** Facial Recognition. Machinery that recognizes identities from facial imagery or videos.
- FS** Father-son. Relationship type (pairwise).
- GFGD** Grandfather - granddaughter. Relationship type (pairwise).
- GFGS** Grandfather - grandson. Relationship type (pairwise).
- GGFGGD** Great GFGD. Relationship type (pairwise).
- GGFGGS** Great GFGS. Relationship type (pairwise).

**GMGD** Grandmother - granddaughter. Relationship type (pairwise).

**GMGS** Grandmother - grandson. Relationship type (pairwise).

**GGMGGD** Great GMGD. Relationship type (pairwise).

**GGMGGS** Great GMGS. Relationship type (pairwise).

**MAP** Mean Average Precision. The average AP across various instances.

**MD** Mother-daughter. Relationship type (pairwise).

**MID** Member ID. Unique identifier assigned to each family member of FIW.

**MS** Mother-son. Relationship type (pairwise).

**PID** Picture ID. Unique identifier assigned to each photo of FIW.

**RFIW** Recognizing FIW. Annual data challenge for recognizing kinship in visual media.

**ROC** Receiver operating characteristic [curve].

**SAM** Soft-argmax. 2D softmax function.

**SIBS** Brother-sister. Relationship type (pairwise).

**SS** Sister-sister. Relationship type (pairwise).

**TAR** True-acceptance rate.

**VID** Video ID. Unique identifier assigned to each video of FIW-MM.

# Acknowledgments

With the highest regards, and the upmost appreciation, are the incredible people in my life. My life has been blessed with diverse network of brilliant, passionate, and sincere people from all over the globe. A sub-population of which are also a part of my professional network; with many more personal. Let me now take a moment for me to express the deep appreciation I have for all of those that have had an impact on my life as far as the pages of this thesis goes. Hence, I owe many thanks to many that helped with completion of my dissertation, for they had involvement that directly relates to the success of this thesis. Whether they realized it or not, I would not be in the position I currently find myself as I prepare this report in exchange for a PhD in Computer Engineering.

## MENTORS

During this lifetime, there have been many who inspired, influenced, and impressed me. However, few, if any, come close to matching the level of impact on me as my PhD advisor, Dr. Yun Raymond Fu. For starters, he convinced me to get a PhD, helped get me admitted and registered past deadlines (*i.e.*, allowed since an undergraduate husky with the support of my advisor to be). Now nearly five-years later, and with the knowledge acquired from this decision, there succeeded an extent never imagined. That was hope for a prospective graduate student to gain research experience from Raymond, renowned in research communities and proven very clearly why. Others and I learn under his remarkable insights and knowledge of research. He enables us at high and low-levels: the high-level being the motivation of research, while the low-level is the technical novelty in the mathematical models and algorithms. Nearly every week for the past 5 years, Raymond set his sight in the lab on each PhD (*i.e.*, an average of over ten students at any given time) allowing us to learn from each other and provide time to review progress. Raymond's expertise generally just knows the best route to take (*e.g.*, the best conference to target for paper submission, the proper way to reach out to other research groups to inquire about a need of ours, even who we should try to speak with at an upcoming conference). At the same time, Raymond trusts us, so he often reaches to one of his students for the most accurate responses and feedback. A big emphasis is put on humility and aggressiveness such as teaching us to work for what we want, appreciate it when there, but readily move onward to the next task. Beyond research, Raymond remains an incredible mentor. Ask anyone in our department— Raymond is a tough advisor with a high expectation. Over the years, I have come to realize the amount of extra effort that is for him (*i.e.*, it is much easier and quicker to say little, but Raymond very rarely cuts it cheaply for us). Week-in, week-out, Raymond devotes great efforts into us: pushing us to be the best possible by recognizing our strengths and weaknesses to help us leverage one while improving the other. He is a trustworthy mentor, and a great friend. I have learned

so many different facets from him (*i.e.*, research, professional, social). I witnessed countless alumni go through the process under his continuous guidance, advice, effort, patience, and encouragement.

I would also like to thank the rest of my PhD dissertation committee, Professor Sarah Ostadabbas and Professor Octavia Camps. Both provided constant support in preparing and delivering my dissertation, while getting opportunities to get to know them through independent means. Specifically, the first time I was exposed to a computer vision topic it was per one of Professor Camp's courses during undergrad. From the get go, Professor Camps was completely approachable. Furthermore, her knowledge in machine vision had inspired me to push myself to learn all that is possible. Along with several vision and related course, Professor Camps also oversaw the oral portion of my qualifying exam completed second year of graduate. Furthermore, at the start and until today I regularly see Professor Camps at vision conferences— it is always a pleasure to catch up with Professor Camps, whether we are on campus or a venue in South Korea. There is never a time she waves me away: I appreciate deeply her expressing interest in my research and status as a graduate student, for she will always check in and have discussions when we bump in to one another. The same for Professor Ostadabbas: another faculty active at conferences and events. Sara Ostadabbas always shares useful bits of information, whether through formal or informal discussions—for instance, during a conference workshop we both chaired, we were sitting together up front taking turns introducing the next speaker, at which time Professor Ostadabbas kept elaborating on several points she raised during the PhD dissertation proposal presentation I had recently delivered. This really meant a lot: several weeks after my proposal, and Professor Ostadabbas still could recall specifics to further elaborate on advise she had given. Along with a good memory, she clearly cared. Both Professor Camps and Professor Ostadabbas clearly have exceptional amounts of care for us students, and the studies as a whole. The two of them alongside my advisor Raymond made for an all-star PhD committee that really made a difference for me in the end.

Professor Ming Shao, although not a part of my PhD committee, has been an outstanding mentor of mine for many years. I was fortunate that the year I joined the lab was about the year Professor Shao was graduating: when we started working together he was a senior student of SMILE. Thankfully, he remained in academia as a professor at University of Dartmouth. Ming and I have collaborated on many works and efforts: from the first kinship paper on FIW to the most recent challenge. Especially in cases the task is new to me (*e.g.*, first journal rebuttal or organized workshop), Ming has been outstanding at teaching me to understand (*i.e.*, not to just complete this time, but to get every time). Lucky for me: Ming has grown more of a friend than a colleague.

## COLLABORATORS AND COLLEAGUES

I would like to thank all the members of the SMILE Lab whom I have had the pleasure of working directly with, such as Dr. Handong Zhao, Yu Yin, Can Qin, Yulun Zhang, Professor Sheng Li, Professor Zhengming Ding, Lichen Wang, Kunpeng Li, and Zaid Khan— many of whom I had great moments traveling with on behalf of the work. Also, the many other members of SMILE Lab for which I spent endless hours alongside working independently like Dr. Kang Li, Dr. Chengcheng Jia, Professor Hongfu Liu, Dr. Shuyang Wang, Dr. Shuhui Jiang, Dr. Jun Li, Professor Yu Kong, Dr. Zhiqiang Tao, Songyao Jiang, Bin Sun, Professor Yi Tian, Professor Qianqian Wang, Gan Sun, Kunpeng Li, Kai Li, Huixian Zhang, Zhenglun Kong, and Chang Liu. Grouped together by Dr. Fu

into SMILE Lab, this hard-working, highly-expecting individuals, yielded a closely knit, synergistic nature in an inspiring environment and allowed many to successfully prosper.

I was extremely lucky to get matched up with kind and intelligent people bosses on interns-mentors that taught me vast types of knowledge. Acknowledging the most recent, graduate-level advisors: Samson Timenor (ISM Connect), Sergey Tulyakov (Snap Inc.), Jeffrey Byrne (STR). Having interned with ISMConnect as a senior grad student, the time could not have been more perfect for the many lessons learned from Samson: beyond technical, and full of life-long concepts and snippets of wisdom to better shape me for the professional (and personal) life to come. While at Snapchat, Sergey and I spent lots of time together to draft a paper, end-to-end, in the limited time of summer. Sergey had endless lessons and technical critique that helped form my way of thinking as research scientists (*i.e.*, every experiment needs hypothesis, should be well thought out and with consideration to variables under investigation, and pre-notions for the next steps whether a null or alternative hypothesis results). As simple of a concept, and as obvious as it seems when speaking of, us researchers often find ourselves overwhelmed in thoughts and with a medley of ideas that it is not uncommon to windup saturated, where remaining grounded yields higher quality output of our work– the emphasis Sergey put on this, whether he realized it or not, will forever remain at the forefront of my thought process as a professional and, sometimes, even beyond. Finally, Jeffrey at STR took me on as an intern for consecutive summers my earliest years of graduate school. Even with many projects being restricted to me as a student (*i.e.*, government classified), Jeffrey was able to carve out meaningful projects based on technical concepts that are regularly found useful today (*e.g.*, algorithms to sort and search, along with entire topics such as clustering, adversarial ML, advanced code development and API design in C++ and Python, and more). Furthermore, Jeffrey included me in many meetings– the lessons learned here were exceptional: the ways Jeffrey hosted, conducted, and led team meetings in ways that motivated, inspired, and included everyone. Jeffrey’s tendency to build strong teams– the language used with those who are the best in the respective topic at hand; the organization and communication that tied everything together– it was later clear to me how Jeffrey led the top team in many competitive government programs (*e.g.*, JANUS).

Even before the PhD, but still at Northeastern as an undergraduate, there had been countless scenarios that helped shape me as a whole. Too many to accurately pin-point every instance in a quick writing piece. Nonetheless, few had such outstanding contributions to me as a person that lessons and other impressions from them will likely continue to propagate for many more years to come. Specifically, those that had a direct hand in helping me find and establish myself both personally and professionally had me realizing I could not be more fortunate in my network– Dr. Charles DiMarzio (and his wife Sheila and friend Maureen) by welcoming me into his lab, taking on the role as a mentor and friend, while providing guidance to my young professional-self; Kristin Hicks, having been the very first person I had contact with from Northeastern, with a phone-call welcoming me to a summer research program as a visiting community college student, which wound up being the first of many opportunities with Kristin (including graduate school fellowship); Dr. David Kaeli, with his advice that helped me transfer to NEU, and then later land my first co-op at Analogic; Dr. Richard Harris, and just being an incredible role model for me to watch, learn from, and often talk to (there would be nights in his office after we organized a workshop with him on his computer guiding me in more ways than he probably even imagined), and providing me the opportunity to grow as a public speaker. There are so many others that I hesitated in naming any. However, the ones mentioned have been there from start, and are still there today. Just incredible and so fortunate.

## FAMILY AND FRIENDS

My mother has had such profound impact on me as a whole; had I not been blessed enough to have my mother chances of me aiming for the highest sights, while remaining as happy and healthy as possible, would likely not be at the levels for which it is. Whether it be encouraging me when down, challenging me at times of comfort, or teaching me at times of growth, my mother has always been there– as the gym teacher when in kindergarten, the driver for doctor’s appointments well before the age I could drive, or one of the many times moving while furthering my education– my mother has always there, she always showed up, and she has always loved myself and siblings. Words cannot express to what extent having a mother like mine has had on me as a person– the explanation is infinite, while words are discrete, so I would argue it is impossible for me to express the true extent for which this is meant (physically, emotionally, spiritually, and many other facets encompassed in the nature of my very existence). Beyond my mother, Lisa Robinson, taking on the role of *super mom*, she has also been the greatest teacher of mine in the broad subject of *Life*. Not only has she helped me financially when in need, but she was also there in spirit when I needed a boost; on the other hand, she would also be the first to take me down a peg or further challenge me at times of over comfort. It is almost like she indirectly trained me for Dr. Fu, for Dr. Fu would often remind us to remain “humble, but aggressive.” Looking back at the big picture, such a mindset needed to live by these words lie at foundation of many of the lessons and experiences posed by mother.

Alongside my mother were my siblings Stacey (Robinson) McGuire, Thomas Robinson, Brendan Crocker, and Briar Crocker, whom are also amongst my best friends. Also, my fathers Peter Robinson and Paul Crocker, my Aunt Theresa Robinson, my cousins John Robinson and Lorri Robinson, and my grandparents Tom and Patricia Floramo who helped make this dream a reality. Also, my girlfriend through the completion of my dissertation and beyond, Briana, who has been by my side through the entire PhD process: dealt with me being consumed by deadlines, and been there as my support at moments of feeling overwhelmed or burned out.

Many family members and friends have supported me by peer reviewing papers coherency and language– my brother Thomas Robinson and my friends Laura Rose, Maureen Hawe, and Bruce Collins. All of whom taught me indirectly by providing feedback on my writing pieces. Especially towards the end, when Laura Rose went over and above helping me polish up language in the final papers published as a graduate student, including this work (*i.e.*, this thesis, end-to-end).

Study buddies friends who I spent countless hours with helped motivate and push me to higher limits: Jordan Kiellach, Brian Toner, Juan Ramirez, Ryan Snyder, Joshua Mcdougall, Robert Watson, and many, many more. Last, but certainly not least, are those that supported me via letters of recommendation and as references, both personal and professional, for one or more of the many endeavors that led me to today: Bonnie-Jeanne Toner, Richard Harrison, Chris McQuire, Samson Timoner, Yun (Raymond) Fu, David R. Kaeli, Charles A. DiMarzio, Kristin Hicks, Paul Chandley, Clair Duggan, Nancy Nickerson, and certainly others over many years of doing co-ops, research, and other rich experiences made possible as a member of the NEU community. Additionally, are those who helped me move many times between undergraduate and graduate tracks at NEU: Paul Crocker, Lisa Robinson, Bruce Collins, Theresa Robinson, and John Robinson.

I thank you all for your love, patience, and support. I must acknowledge it, for it was by you the seemingly impossible goal of yesterday (*i.e.*, a doctor of engineering) has become the reality of today. I hope to return the favor in years of research to come, *i.e.*, let me apply the many years of training and life experiences to making the world a better place come tomorrow.

# Abstract of the Dissertation

Automatic Face Understanding: Recognizing Families in Photos

by

Joseph P Robinson

"Doctor of Philosophy" in Computer Engineering

Northeastern University, February 11, 2023

Dr. Yun Fu, Advisor

Visual kinship recognition has an abundance of practical uses. For this, we built the largest database for kinship recognition, FIW. Built entirely in-house with no cost using a semi-automatic labeling scheme. Specifically, we first aligned faces detected in family photos with names in the corresponding text metadata to mine the label proposals with high confidence. The remaining data were labeled using a novel clustering algorithm that used label proposals as side information to guide more accurate clusters. Great savings in time and human input was had. Statistically, FIW shows enormous gains over its predecessors. We have several benchmarks in kinship verification, family classification, tri-subject verification, and large-scale search & retrieval. We also trained CNNs on FIW and deployed the model on the renowned KinWild I and II to gain state-of-the-art (SOTA). Most recently, we further augmented FIW with multimedia (MM) for 200 of its 1,000 families- a labeled collection we dubbed FIW-MM. Now, video dynamics, audio, and text captions can be used in the decision making of kinship recognition systems.

FIW continues to pave the way for this research track: (1) advanced SOTA (*e.g.*, marginalized denoising auto-encoder based on metric learning that preserves intrinsic structures of kin-data and encapsulates discriminating information in learned features); (2) introduced generative models to predict a child's appearance from a parent pair (*i.e.*, proposed an adversarial autoencoder conditioned on age and gender to map between facial appearance and these higher-level features for control of age and gender); (3) designed evaluations with benchmarks to support challenges, workshops, and tutorials at top tier conferences (*e.g.*, CVPR, MM, FG, ICME), and a premiere Kaggle Competition. We expect FIW will significantly impact research and reality.

Additionally, we tackled the classic problem of facial landmark localization in images. This is a task that has been in focus for decades, and many solutions have been proposed. However, there are revamped interests in pushing facial landmark detection technologies to handle more challenging

data with deep networks now prevailing throughout machine learning. A majority of these networks have objectives based on L1 or L2 norms, which inherit several disadvantages. First of all, the locations of landmarks are determined from generated heatmaps (*i.e.*, confidence maps) from which predicted landmark locations (*i.e.*, the means) get penalized without accounting for the spread: a high scatter corresponds to low confidence and vice-versa. To address this, we introduced a LaplaceKL objective that penalizes for low confidence. Another issue is a dependency on labeled data, which is expensive to collect and susceptible to error. We addressed both issues by proposing an adversarial training framework that leverages unlabeled data to improve model performance. Our method claims SOTA on renowned benchmarks. Furthermore, our model is robust with a reduced size: 1/8 the number of channels (*i.e.*, 0.0398 MB) is comparable to state-of-the-art in real-time on a CPU. Thus, our method is of high practical value to real-life applications.

Finally, we built the Balanced Faces in the Wild (BFW) dataset to serve as a proxy to measure bias across ethnicity and gender subgroups, allowing us to characterize FR performances per subgroup. We show performances are non-optimal when a single score threshold is used to determine whether sample pairs are genuine or imposter. Furthermore, actual performance ratings vary greatly from the reported across subgroups. Thus, claims of specific error rates only hold for populations matching that of the validation data. We mitigate the imbalanced performances using a novel domain adaptation learning scheme on the facial encodings extracted using SOTA deep nets. Not only does this technique balance performance, but it also boosts the overall performance. A benefit of the proposed is to preserve identity information in facial features while removing demographic knowledge in the lower dimensional features. The removal of demographic knowledge prevents future potential biases from being injected into decision making. Additionally, privacy concerns are satisfied by this removal. We explore why this works qualitatively with hard samples. We also show quantitatively that subgroup classifiers can no longer learn from the encodings mapped by the proposed.

# Chapter 1

## Introduction

When your face says it all, your mouth  
waits its turn.

---

*Anthony T. Hincks.*

As known by many, and often in the form of common sense, facial cues hold an abundance of information— whether it be the identity of the subject, their age, or even the way they feel in the moment. Hence, the human face, as a biometric, holds high potential in its relevance in vast practical use-cases. For starters, automatic face understanding makes up an exceptionally large problem space: face-based problems can be as specific as identity classification, or no interest in the identity but more broadly as it as an object (*i.e.*, face detection). Now, beyond the more traditional problems of identification and detection exist a slew of attribute-based tasks. As mentioned, lots can be learned from facial cues, which range from the measure of the presence or absence of an emotion, comfort or pain, honest or adversary, focused or distracted, and even rested or exhausted. Furthermore, facial cues (*i.e.*, faces captured in imagery or video data) encapsulate knowledge of intrinsic characteristics or attributes, such as gender and age; extrinsically faces can relate to one another by grouping. For instance, determine whether two or more subjects are blood relatives by comparing face data.

Nowadays, many use-cases for face data in machine vision have been explored— face-based technology can be seen throughout society, and in various forms. To name a few: used to unlock mobile device, provide access control in a security sensitive setting, and automatically tag on a social media platform. For any of these to have been possible, and to match the high capacity set by the modern-day, data-driven deep learning models, many have spent effort and resources to acquire and share large-scale benchmark facial image collections. Hence, some of the great advances in

## CHAPTER 1. INTRODUCTION

automatic FR technology would never have been achieved without the large-scale datasets of labeled faces: of the many face-based problems there are opportunities in research and venture capital that can then be achieved. The same holds for the face-based attributes: due to the fine-grained nature of face data, along with a large number of samples readily available to scrape from the web, there exists a large variance inherent naturally by the data. As we will discuss in the later chapters on bias, consideration for the large variance, along with controlled variables that allow us to minimize bias, is critical. Thus, the same factor that motivates us upfront relates to a problem that stems up thereafter (*i.e.*, the need of big data and the effects of working with big data in face-based problems).

When drafting up the research questions (or topic), we reviewed the state of various face-based problems. When drafting a research plan the winter of 2015, kinship recognition from faces, in particular, caught my attention— automatic FR of nearly every type is used in main-stream solutions. For instance, identification system for recommending tags on social media platforms, or a face tracking and landmark detection system for apps like Snapchat. A FR-based view that received minimal attention was kinship recognition. A problem statement we next will define.

### 1.1 Objective and Significance

Our goal was to acquire a framework to detect pairs of kinship from facial images. Furthermore, we intend to bridge the gap separating research-from-reality by working to develop machinery to automatically detect kinship via face data in visual media. Put differently: we set out to develop a system that accepts two or more face images as input, and outputs the class of KIN if blood relatives and NON-KIN otherwise. We focus on an analysis that highlights different cases brought on by the evaluation metrics discussed as a part of Part I (Chapter 2). As part of this dissertation, we aim to determine how well kinship can be recognized from facial cues. As we will see, a richer resource (*i.e.*, paired data) was required to sufficiently model kinship from sets of faces from different subjects.

#### 1.1.1 Scope

There are various views of the problem space for visual kinship recognition. In essence, the different views are organized using different data splits, label types, and metrics. Typically, the problem was either crafted to fit the data, or the data was shaped for the problem statement. Also, existing settings of related tasks for similar problems are often borrowed or used for inspiration (*e.g.*, LFW identification benchmark inspired data splits for our FIW image collection).

## 1.2 Contributions

Most of this dissertation covers work previously published in peer reviewed journals and conference proceedings. For convenience, all work published as a part of this effort are listed at the end of this section. Besides in the following subsection, citations used throughout the report are in reference to the main bibliography at the end of the document– the list of papers provided here is for quick reference as we review contributions and describe organization. Again, it is important to note that this section is the only part of the manuscript that uses the paper list at the end of this section.

### 1.2.1 Organization

The dissertation consists of three parts: *Preliminary*, *Processing*, and *Post-processing*. Let us now explain the contents of each. Note that references in the following paragraph refer to the list provided in the previous section. This is the last time the list of personal publications is referred to, and the bibliography contains the only citations used from here onward.

**Preliminary (Part I).** We start by reviewing FR as a whole in Chapter 2. There are many topics common between traditional and kinship-based face recognition, all of which we introduce first and foremost. This leads to the pre-processing of faces (*i.e.*, face detection and alignment) that is required at the beginning of a FR system (Chapter 3). For the discussion we hone-in on our work in [11].

**Process (Part II).** Covers kinship recognition. Specifically, we review existing work in visual kinship recognition (Chapter 4), and then our contributions with FIW dataset and SOTA benchmarks. A number of our works fall in (Chapter 5), *i.e.*, [2], [12], [13], [14], [15], [16], [17], [18]. Then, in Chapter 6, we discuss the recent release of the labeled multimedia for subjects of FIW: FIW-MM, along with the semi-automatic machinery used to label the multi-modal data with minimum human inputs and no financial costs [3].

**Post-processing (Part III).** We do qualitative studies on the different results in kinship recognition in Chapter 7. This naturally leads to a discussion on the limitations of SOTA, along with the technical challenges currently at the forefront. The aforementioned were findings from [2], [7], [14].

Finally, our studies on bias in FR systems is reviewed (Chapter 8). For this, we built and shared the novel Balanced Faces in the Wild (BFW) face dataset to benchmark facial recognition systems with balanced data [8]. In addition to the machinery being bias, we also measure bias inherent in humans as well– we conducted a human survey to measure and analyze the human

## CHAPTER 1. INTRODUCTION

bias. Furthermore, we describe a novel feature adaptation technique we proposed to mitigate issues from unbalanced performances across subgroups [1]. Additionally, our debiasing technique also benefits in areas of privacy and protection– the objective used to adapt features involves a penalty for recognizing the subgroup. Thus, the resulting mapping of the debiases and the features removes knowledge of the particular subgroup as well.

We conclude with a discussion on next steps while concluding the various works represented and discussed in this thesis (Chapter 9).

### 1.2.2 Publications

Publications are listed in reverse-chronological order. Each item available online has ‘[paper]’ appended, which provides a direct link to the respective paper along with references to the main bibliography.

1. **Joseph P. Robinson**, Can Qin, Yann Henon, Samson Timoner, and Yun Fu. “Balancing Biases and Preserving Privacy on Balanced Faces in the Wild,” in *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2020 (Under review).
2. **Joseph P. Robinson**, Ming Shao, and Yun Fu. “Visual Kinship Recognition: A Decade in the Making,” CoRR arXiv:2006.16033, 2020. (Under review, Trans. on PAMI). [\[paper\]](#) [\[18\]](#)
3. **Joseph P. Robinson**, Z. Khan, Y. Yin, M. Shao, and Yun Fu, “Families in wild multimedia (FIW-MM): A multi-modal database for recognizing kinship,” *CoRR arXiv:2007.14509*, 2020. (Under review, Trans. on MM). [\[paper\]](#) [\[19\]](#)
4. Yu Yin, **Joseph P. Robinson**, and Yun Fu, “Multimodal In-bed Pose and Shape Estimation Under the Blankets,” CoRR arXiv:2012.06735, 2020. [\[paper\]](#) [\[20\]](#)
5. Yu Yin, **Joseph P. Robinson**, Songyao Jiang, Yue Bai, Qin Can, and Yun Fu, “SuperFront: From Low-resolution to High-resolution Frontal Face Synthesis,” CoRR arXiv:2012.04111, 2020. [\[paper\]](#) [\[21\]](#)
6. Chengyao Zheng, Siyu Xia, **Joseph P. Robinson**, Changsheng Lu, Wayne Wu, Chen Qian, and Ming Shao. “Localin Reshuffle Net: Toward Naturally and Efficiently Facial Image Blending,” in 15-th Asian Conference on Computer Vision (ACCV), 2020. [\[paper\]](#) [\[22\]](#)

## CHAPTER 1. INTRODUCTION

7. Yu Yin, Songyao Jiang, **Joseph P. Robinson**, and Yun Fu “Dual-attention GAN for large-pose face frontalization,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2020. [[paper](#)] [[23](#)]
8. Lichen Wang, Bin Sun, **Joseph P. Robinson**, T. Jing, and Yun Fu, “Ev-action: Electromyography-vision multi-modal action dataset,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2020. [[paper](#)] [[24](#)]
9. **Joseph P. Robinson**, Yu Yin, Zaid Khan, Ming Shao, Siyu Xia, Michael Stopa, Samson Timoner, Matthew A. Turk, Rama Chellappa, and Yun Fu, “Recognizing Families In the Wild (RFIW): The 4th Edition,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2020. [[paper](#)] [[7](#)]
10. **Joseph P. Robinson**, G. Livitz, Y. Henon, C. Qin, Yun Fu, and S. Timoner, “Face recognition: too bias, or not too bias?” in *Conference on Computer Vision and Pattern Recognition Workshop*, 2020. [[paper](#)] [[25](#)]
11. Yu Yin, **Joseph P. Robinson**, Zhang, Y., and Yun Fu. “Joint super-resolution and alignment of tiny faces,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020. [[paper](#)] [[26](#)]
12. W. Zhuang, Y. Wang, **Joseph P. Robinson**, C. Wang, M. Shao, Y. Fu, S. Xia. “Towards 3D Dance Motion Synthesis and Control,” in *CoRR arXiv preprint arXiv:2006.05743*, 2020. [[paper](#)] [[27](#)]
13. **Joseph P. Robinson**, Yuncheng Li, Ning Zhang, Yun Fu, and Sergey Tulyakov, “Laplace landmark localization,” in *IEEE International Conference on Computer Vision (ICCV)*, 2019. [[paper](#), [poster](#)] [[28](#)]
14. P. Gao, S. Xia, **Joseph P. Robinson**, J. Zhang, C. Xia, M. Shao, and Yun Fu, “What will your child look like? DNA-net: Age and gender aware kin face synthesizer,” *CoRR arXiv:1911.07014*, 2019. [[paper](#)] [[10](#)]
15. Yue Wu, Z. Ding, H. Liu, **Joseph P. Robinson**, and Yun Fu, “Kinship classification through latent adaptive subspace,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2018. [[paper](#)] [[29](#)]
16. **Joseph P. Robinson**, Ming Shao, and Yun Fu, “To recognize families in the wild: A machine vision tutorial,” in *ACM Conference on Multimedia (ACMMM)*, 2018. [[paper](#)] [[30](#)]

## CHAPTER 1. INTRODUCTION

17. **Joseph P. Robinson**, Ming Shao, Yue Wu, Hongfu Liu, Timothy Gillis, and Yun Fu, “Visual kinship recognition of families in the wild,” in *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2018. [[paper](#)] [6]
18. **Joseph P. Robinson**, M. Shao, H. Zhao, Y. Wu, T. Gillis, and Yun Fu, “Recognizing families in the wild (RFIW),” in *RFIW at ACM MM*, 2017. [[paper](#)] [5]
19. S. Wang, **Joseph P. Robinson**, and Yun Fu, “Kinship verification on families in the wild with marginalized denoising metric learning,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2017. [[paper](#)] [31]
20. **Joseph P. Robinson**, M. Shao, Y. Wu, and Yun Fu, “Families in the wild (FIW): Large-scale kinship image database and benchmarks,” in *ACM Conference on Multimedia (ACMMM)*, 2016. [[paper](#)] [5]
21. **Joseph P. Robinson**, Edward Scott, and Yun Fu, “Pre-trained D-CNN Models for Detecting Complex Events in Unconstrained Videos” in *SPIE Commercial and Scientific Sensing & Imaging*, 2016. [[paper](#)] [32]
22. **Joseph P. Robinson**, M. Shao, Y. Wu, and Yun Fu, “NEU- MITLL @ TRECVID 2015: Multimedia Event Detection by Deep Feature Learning,” in *Proceedings of TRECVID, NIST, USA*, 2015. [[paper](#)] [33]

### 1.2.3 Service

Here are selective services completed as part of this dissertation. Having served as a reviewer on several journals (*e.g.*, IEEE Trans. on PAMI, TIP, many others), PC or SPC for many conferences (*e.g.*, CVPR, AAAI, IJCAI, ICCV, ECCV, FG (3x awarded *Outstanding Reviewer*, and many others), and for many years, these listings were omitted and are, thus, not explicitly listed. Instead, workshops, challenges, and tutorials at top tier conferences, for which my contributions were critical to the success of the event are listed as follows.

#### 1.2.3.1 Workshops

1. 2021 Workshop Chair, *10th Workshop on the Analysis & Modeling Faces & Gestures (AMFG)*, CVPR (online). [[web](#)]

## CHAPTER 1. INTRODUCTION

2. 2020 Challenge Chair and Organizer, *4th Recognize Families In Wild (RFIW) Challenge*, IEEE FG Argentina. [[web](#)]
3. 2019 Host & Organizer, *Recognizing Families in the Wild*, CVPR Long Island, CA. [[web](#)]
4. 2019 Tutorial Host & Organizer, *Recognize Families: A Machine Vision Tutorial (II)*, IEEE FG Lille, France. [[web](#)]
5. 2019 Workshop Chair, *9th Workshop on AMFG*, CVPR Long Island, CA. [[web](#)]
6. 2019 Workshop Chair, *2nd Workshop on Faces in Multimedia (FacesMM)*, ICME Shanghai, China. [[web](#)]
7. 2019 Challenge & Workshop Chair, *3rd RFIW Challenge*, IEEE FG Lille, France. [[web](#)]
8. 2018 Host & Organizer, *RFIW: Machine Vision Tutorial*, ACM MM Seoul, S. Korea. [[30](#)]
9. 2018 Host & Organizing Chair, *1st Workshop on FacesMM*, ICME San Diego, CA. [[web](#)]
10. 2019 Challenge & Workshop Chair, *2nd RFIW Challenge*, IEEE FG China. [[web](#)]
11. 2018 Workshop Host & Organizing Chair, *8th Workshop on AMFG*, CVPR SLC, Utah. [[web](#)]
12. 2017 Host & Organizing Chair, *New England CV Workshop*, NEU Boston, MA. [[web](#)]
13. 2017 Host & Organizer, *RFIW Data Challenge Workshop*, ACM MM Mountain View. [[web](#)]

## **Part I**

# **Preliminaries**

## Chapter 2

# Automatic Facial Recognition

### 2.1 Overview

To describe our contributions in automatic face recognition and understanding technology, an overview of fundamentals, such as problem statements, related face-based systems and efforts, along with data preparation and evaluation concepts that relate to the work done for this dissertation.

We start by reviewing major milestones in conventional FR. A basic understanding in concepts pertaining to conventional FR is imperative for understanding kinship recognition from image data (*i.e.*, facial images). For this, a brief look at traditional systems before those that are more popular nowadays. As we will discuss, deep learning-based approaches are mostly data-driven: feature learning, a large model complexity that demands more data to avoid over-fitting (*e.g.*, Convolutional Neural Networks (CNNs), generative adversarial networks (GANs), and much much more). As a part of the basic depiction of the aforementioned is a section of the different loss functions that had previously claimed SOTA and had later been employed in a face-based kinship recognition system. Performance ratings for both the traditional and the modern-day FR methods are reported for face identification task to provide insight in our expectations later when used in kinship recognition– this is especially true for loss functions, which has been pivotal in FR.

We then step back a couple of modules in the ML pipeline to preprocessing. Most face-based systems depend on a face detector (*i.e.*, object detector, with face as a boolean target for whether it is present or not). Furthermore, specific landmarks of the face are too detected: facial landmarks are often used to crop brother-brothers (BBs) more consistently across a large set; also, landmarks are treated as points of reference to frontalize (*i.e.*, align) faces in images prior to feeding to model. Considering the preprocessing of faces (*i.e.*, detection, aligning, etc.). is conducted on

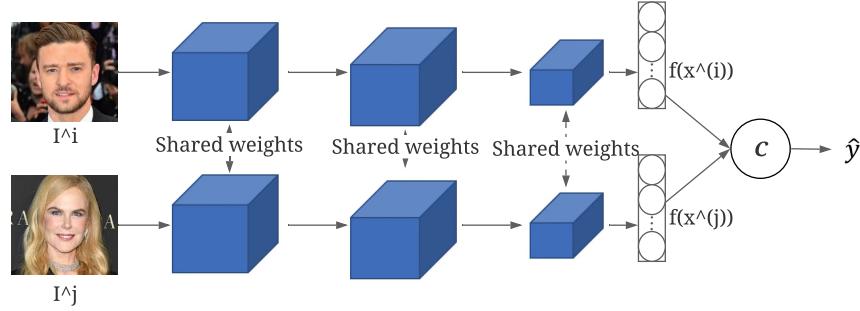


Figure 2.1: **A generic verification system.** Regardless of the network specifications (*i.e.*, independent of layer counts, layer sizes, type of metric set topmost). The aim is to map two image inputs to a single logistic value. Input images (*i.e.*,  $I^i$  and  $I^j$ ) must then assume the same size as that of the input layer of the model. Then, the learned representation of the  $i - th$  and  $j - th$  faces are output from the last layer of the network prior to the classification layer. An aspect of the model shown is the mechanism used to fuse (*i.e.*, single logit value output from  $C$ , which was fed the features of the  $i - th$  and  $j - th$  face encoding). The framework is inherently boolean, as the task is to map a sample pair (*i.e.*, pair of faces) to a boolean class label (*i.e.*, 1 if *genuine* and 0 for *impostor*).

nearly all of our work, while a majority of all others (*i.e.*, for conventional face recognition, along with nearly all other discriminative tasks like kinship recognition using facial cues). Finally, we built several large-scale face imagesets— an effective preprocessing setup shows to be exceptionally important in problem spaces of increasing difficulty.

All topics covered in this chapter are meant to be span out broadly. In other words, this preliminary information is not the core (*i.e.*, meat) of this dissertation. However, the basic understanding we hope to provide the reader is believed to be essential for the topics covered in the later chapters. Furthermore, several efforts spent on this dissertation fall in the realm of face-based preprocessing. On the one hand, our papers that relate to the topics of this chapter are included. Still, the thesis is not on these works and will not tailor the upcoming discussion.

Having recently joined the researchers of the Machine Vision community, I have already been exposed to a wide range of contemporary problems tackled by both classic and modern methods. However, up to this past month I have had no experience working with faces. That was until a recent research assignment involving kinship recognition from digital photos came about. Hence, understanding faces in the eyes of computer vision is essential here. Plus, as a young researcher in the field, I felt it was essential to gain some understanding of face detection and facial recognition.

## CHAPTER 2. AUTOMATIC FACIAL RECOGNITION

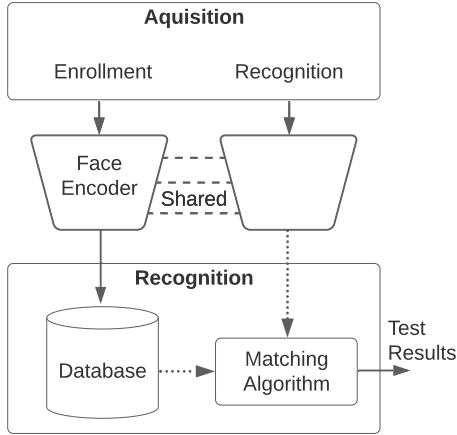


Figure 2.2: **Generic recognition system.** Faces are often stored as encodings in a database through *enrollment*. *Recognition* is then to compare the encoding of an input face to those in the database.

Due to both space and time constraints and, moreover, relevance to the contributions of the dissertation, only a few FR et *Eigenfaces*, *Fisherfaces*, Local Binary Pattern Histograms. Modern Facial Recognition The goal of human facial recognition is to automatically verify or identify an individual from digital data, *i.e.*, images or videos (image stacks). Facial recognition, in itself, is applicable in a wide range of technologies. Much of this applies to security-based applications (*e.g.*, biometric identification systems, human tracking, and surveillance systems of all sorts). Nonetheless, its uses span well beyond the scope of just security-related problems. For example, Facebook uses face verification to suggest its users when a photo is uploaded. In addition, facial verification is largely used in search and retrieval-based tasks, involving images, videos, or even both. Facial recognition is also directly applicable in applications involving kinship (*e.g.*, sorting a family photo album, determining ethnicity, and such) [34]. A purpose for facial recognition is common in areas involving facial images, which have and continue to increase rapidly in this "mobile age" [35].

Common challenges faced with facial recognition can be generalized into two groups, external disturbances (*e.g.*, changes in illumination, occlusions from glasses, beards, makeup, etc.) and internal influences (*e.g.*, facial expressions, head rotations, human aging, etc.). It are these challenges that decipher the pros and cons of any given facial recognition approach, as will be exemplified through analysis of the methods covered in this report.

The following sections introduce three modern approaches to facial recognition, which are *Eigenfaces*, *Fisherfaces* and LBPH. Following this, experimental results running these methods on two face databases are shared and further analyzed.

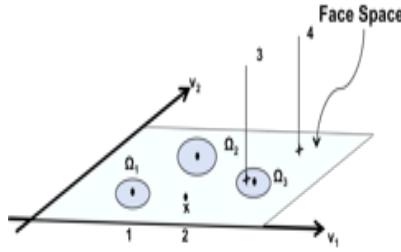


Figure 2.3: **Illustration of the Face Space, which spans the area of the light-blue rectangle.** Note that the blue spheres in the Face Space ( $\Omega_1, \Omega_2, \Omega_3$ ) describe a face of a particular person, *i.e.*, used to identify or verify individual instances. Beyond that, but within the bounds of the Face Space, are varying faces of unknown type, *i.e.*, could be used, for instance, to detect any face.

## 2.2 Traditional Methods

### 2.2.1 Eigenfaces

#### 2.2.1.1 Overview

*Eigenfaces*: the eigenvectors of face images introduced in the early 1990s in research as a FR system. *Eigenface*-related concepts are still widely used today. It is a simple scheme that drastically reduces the search space. Moreover, it performs reasonably well with non-complex face data— a fair assumption to make in many real-world applications based on facial images.

In short, *Eigenfaces* was motivated by the idea that highly dimensional face images are highly correlated, from which the size of the image space can be projected to a space that ignores common features, *i.e.*, preserve the feature dimensions that have less correlation, resulting in a higher variance. In essence, *Eigenface* are based on a dimensionality reduction that selects the features of greatest variance from all features. In turn, preserve fewer features with most information.

Principle Component Analysis (PCA), a dimensionality reduction technique, finds the subspace, projects a feature encoding of an image, and compares to other images via the Euclidean distance or another similarity or distance metric. In other words, the difference between *Eigenfaces* is found and compared to a threshold to determine a final decision, match or no match. This PCA-based dimensionality reduction scheme, generally speaking, uses lower-dimensional vectors to represent some high-dimensional data. In our case, the high-dimensional data are the original images, and the low-dimensional vectors are the *Eigenfaces*. These feature vectors are eigenvectors and the face images projected into this lower-dimensional face-space have been appropriately named *Eigenfaces*,



Figure 2.4: **By taking the difference between each face and the mean of all faces *Eigenfaces* are normalized.** Results in images of facial structure, and resembling a ghost—some call these Ghostfaces, it is most common, and agreeably more appropriate, to refer to them as *Eigenfaces*.



Figure 2.5: **Eigenfaces for the same face under different lighting conditions.** The original face images and their *Eigenface* equivalent are shown on top row and bottom row, respectively. Notice the variation between *Eigenfaces*.

*i.e.*, given a collection of facial images, PCA returns a set of corresponding basis vectors that are used to describe the faces in a lower-dimensional subspace.

In summary, given a face image of size  $N \times N$  (pixels) in its vector exists as a single point in  $N^2$ -dimensional image space. Considering that facial images contain correlated features, these projections shall not be distributed in the large image space, but yet a smaller, greatly reduced space—Such to capture the variations and remove the redundancy of the set of facial images.

In essence, it is a process of minimizing the non-diagonal elements of the covariance matrix, while maximizing the diagonal, which results in a form similar to the following:

$$W = PX = Cov(W) = \frac{1}{n-1} WW^T = \begin{bmatrix} * & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & * \end{bmatrix}, \quad (2.1)$$

where the rows of  $P$  are the principal components of data  $X$ , which preserves most of the information of  $X$ , *i.e.*,  $W$  is the diagonal, also known as the *Eigenfaces*.

### 2.2.1.2 Learning face space

*Eigenfaces*— or as some call it, Ghostfaces (Figure 2.4)—maps a facial image to its *Eigenface* as follows:

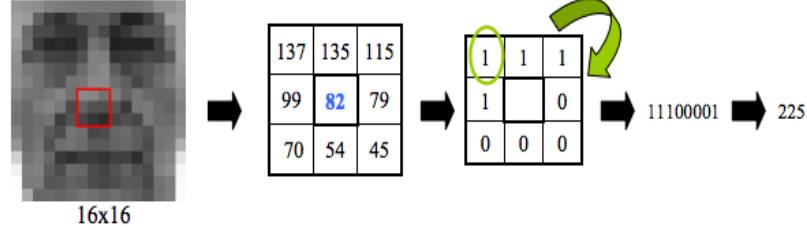


Figure 2.6: Illustration depicting the process of obtaining LBH features from an image for a single pixel. With every pixel processed, a histogram is generated to represent each image.

1. Collect/ obtain a set of  $M$  face images to use for training,  $I_1, I_2, \dots, I_M$ .
2. Vectorize each  $N \times N$  image ( $I_i$ ) to span  $N^2 \times 1$  dimensional space ( $\Gamma_i$ ).<sup>1</sup>
3. Calculate the average face ( $\Psi$ ) of all  $M$  images.

$$\Psi = \frac{1}{M} \sum_{i=1}^M \Gamma_i. \quad (2.2)$$

4. Subtract the mean face from all facial images, *i.e.*, normalize feature vector.

$$\Phi_i = \Gamma_i - \Psi. \quad (2.3)$$

5. Find the covariance matrix of all mean-shifted facial vectors.

$$S = \frac{1}{M} \sum_{n=1}^M \phi_n \phi_n^T = \Phi \Phi^T, \quad (2.4)$$

where  $\Phi = [\phi_1, \phi_2, \dots, \phi_M]$  is a complete set of orthonormal eigenvectors spanning  $N^2 \times M$  and  $S$  is size  $N^2 \times N^2$ .

The mean-shifted image features of each is a linear combination of the eigenvectors, *i.e.*,

$$\hat{\Phi} - \Psi = \omega_1 \mathbf{u}_1 + \omega_2 \mathbf{u}_2 + \dots + \omega_{N^2} \mathbf{u}_{N^2}. \quad (2.5)$$

---

<sup>1</sup>Note that all images are required to have the same resolution, *i.e.*, a preliminary step would be to resize all images to the same size.

Each face can be represented using just the top  $K$  eigenvectors, reducing the dimension of the problem from  $N^2$  to  $K$ , where  $k \ll N^2$  and, hence, is approximated as follows:

$$\hat{\Phi} - \Psi = \omega_1 \mathbf{u}_1 + \omega_2 \mathbf{u}_2 + \cdots + \omega_{N^2} \mathbf{u}_{N^2}, \quad (2.6)$$

or equivalently  $\hat{\Phi} - \Psi = \sum_{j=1}^K \omega_j \mathbf{u}_j$ , and where the  $j^{th}$   $u$  and  $\omega$  is the *Eigenface* and eigenvalue, respectively: the  $k$  eigenvector corresponding to the largest eigenvalues are kept.

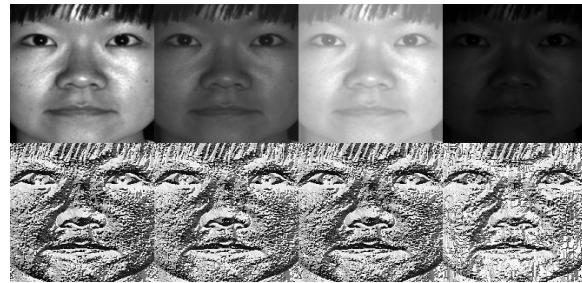


Figure 2.7: **LBPH of the same face under different lighting conditions.** Notice the light invariance that is inherited with this feature.

### 2.2.1.3 Facial Recognition

Provided an unseen image, and a subspace  $\Omega$  that was found during training, images are projected to the Face Space by the formula:

$$Y = \Omega^T (X - \Psi),$$

where

$$\Omega_i = \begin{bmatrix} \omega_1 \\ \omega_2 \\ \vdots \\ \omega_K \end{bmatrix}, \text{ for } i \in [1, 2, \dots, M].$$

The distance between  $y$  and each face class is found as the Euclidean distance.

$$\mathcal{E}_k^2 = \|y - y_k\|^2,$$

where  $k = 1, 2, \dots, M$ .

## CHAPTER 2. AUTOMATIC FACIAL RECOGNITION

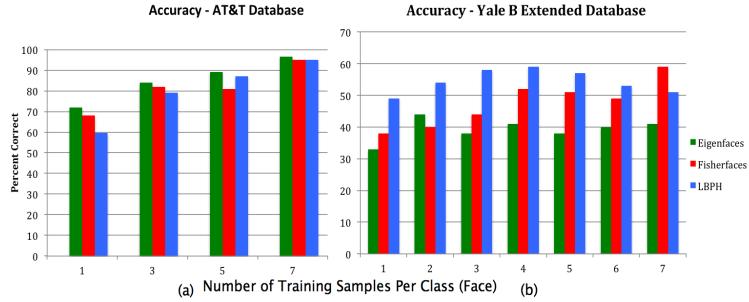


Figure 2.8: **Accuracy measure.** Correct (%) as function of training samples count per class for the AT&T dataset (a). Correct (%) as function of training samples count per class for the Yale B Extended dataset (b).

A distance threshold  $\Theta_c$ , is half the largest distance between any two face images:

$$\Theta_c = \frac{1}{2} \max_{j,k},$$

where  $j = 1, 2, \dots, M$  and  $k = 1, 2, \dots, M$ .

The distance  $\epsilon$  is found between the original image  $x$  and its reconstructed image *Eigenface* from Face space  $x_f$ :

$$\epsilon^2 = \|x - x_f\|^2,$$

where  $x_f = Wx + \mu$ .

### Recognition summarized.

- IF  $\epsilon \geq \Theta_c \Rightarrow$  input image is not a face image;
- IF  $\epsilon < \Theta_c \& \epsilon_k \geq \Theta_c$  for all  $k \Rightarrow$  input image contains an unknown face;
- IF  $\epsilon < \Theta_c \& \epsilon_k^* = \min_k(\epsilon_k) < \Theta \Rightarrow$  input image contains the face of individual  $k*$ .

The main idea is that the  $K$  dimensions embody the most variant aspects across the face images. By this approach, the size reduction of search space greatly outweighs the amount of information compromised. Hence, enough information is preserved in *Eigenfaces* to reconstruct any of the faces used to find the mean-face by doing the above steps in reverse, with, of course, the same mean face that was used to initially shift (normalize) the facial images.

Although *Eigenfaces* work well, there are drawbacks. Considering that we are compressing the information according to the  $K$  eigenvectors with largest variations, *Eigenfaces* are prone to capturing variations caused by external distortions, making it non-invariant to various external

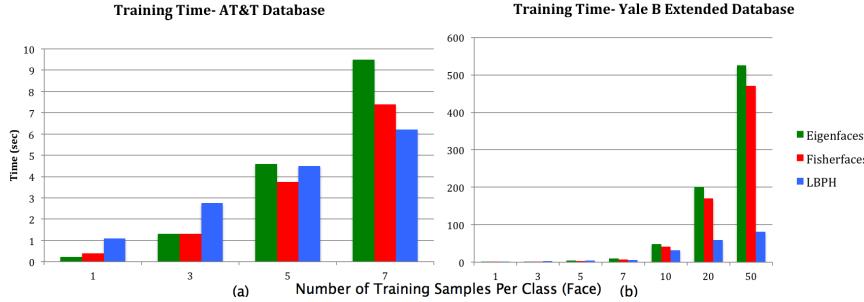


Figure 2.9: **Performance measure.** Training time as function of training samples count per class for the AT&T dataset (Left). Training time as function of training samples count per class for the Yale B dataset (Right).

sources. An example of this shown in Figure 2.5, as the same face in image space is drastically changed in Face Space by changes in lighting.

The following section introduces a method that addressed the drawbacks of *Eigenfaces*—An approach called *Fisherfaces*.

### 2.2.2 Fisherfaces

*Fisherfaces*, like *Eigenfaces*, seek to find a low-dimensional subspace to facilitate the facial recognition task from; *Fisherfaces*, unlike *Eigenfaces*, search for a subspace in a discriminant manner. It overtakes the shortcomings of *Eigenfaces*, while its simple principles are still preserved. It does so by taking the Linear Discriminant Analysis approach [36].

#### 2.2.2.1 Linear Discriminant Analysis (LDA)

LDA finds a basis for projection such that the intra-class variation is minimized, while the inter-class variation is maximized. Rather than explicitly modeling its deviation, we linearly project the high-dimensional data into a subspace that discounts those regions of the face with large deviation. When this approach is used for facial recognition, we refer to it as using *Fisherfaces*. Unlike *Eigenfaces*, and due to the nature of LDA, *Fisherfaces* contain class-information that can be used for classification of the images provided.

LDA class-specific dimensionality approach is summarized as follows:

- Obtain features that best separate between classes, opposed to global scatter.
- Clusters the same classes tightly; Separates between different classes.

- Reduces search space size (*e.g.*, *Eigenfaces*).
- Immune to external agents (unlike *Eigenfaces*).

Given random vector  $\mathbf{X}$  made up of samples from  $c$ -classes.

$$\mathbf{X} = X_1, X_2, \dots, X_c; X_i = x_1, x_2, \dots, x_n.$$

### 2.2.2.2 Within (W)/ Between (B) class scatters

Between-class covaraince:

$$SS_B = \sum_{i=1}^c N_i(\mu_i - \mu)(\mu_i - \mu)^T,$$

where  $\mu = \frac{i}{n} \sum_{i=1}^n x_i$ .

Within-class covaraince:

$$SS_W = \sum_{i=1}^c S_i,$$

where  $S_i = \sum_{x \in D}(x - m_i)(x - m_i)^T$  and  $m$  coordinates retained.

Searches for  $W$  that maximizes class separability criterion:

$$W_{opt} = \arg \max_{\omega} \frac{|W_T S_B W|}{|W_T S_W W|}.$$

### 2.2.3 Local Binary Pattern Histograms (LBPH)

- Uses local (texture) descriptors.
- Thresholds each pixel with neighboring pixels to generate a binary pattern.
- Bins all features for same class into a histogram to obtain a single representation.

The process of obtaining Local Binary Patterns Histograms (PCA) involves a simple calculation, one involving the summation of the differences between the central pixel and all of its neighbors [37].

Unlike *Eigenfaces*, PCA is invariant to changes in lighting (Fig 2.6). The figure shows four images of the same face under different lighting conditions (top row) and the corresponding PCA representation (bottom row). Notice the histograms are nearly identical, regardless of the large variations in light.

### 2.2.4 Results and analysis

#### 2.2.4.1 Experimental Setup

1. Implemented using many MATLAB built-in functions.
2. Built program with two modes:
  - (a) Performance Mode: Experiment varies number of training & test samples.
  - (b) Test Mode: User specifies algorithm, dataset, and number of training samples.
3. Time is measured (seconds) for all calls to training and testing functions.
4. Accuracy measures are based on the number of testing samples.

#### 2.2.4.2 AT&T Database

AT&T is a simple dataset containing mild variations in facial expressions, small rotational shifts, and minor facial obstructions (*e.g.*, glasses). There are a total of 40 subjects, each with 10 face images. The images are 8bit (gray scale) with a size of  $92 \times 112$  pixel, formatted as PGM image files. For every experiment the number of training images is specified, and the remaining (of the 10) are used for testing.

#### 2.2.4.3 Extended Yale B Face Database

Yale B is a more Complex dataset that contains high variation in facial expressions, image rotations, facial obstacles (*e.g.*, glasses and beard), and changes in illumination. There are a total of 38 subjects, each with 64 face images. The images are 8bit (gray scale),  $168 \times 192$ , and PGM formatted. And again, for each experiment the number of training images is specified, and the remaining (of the 10) are used for testing.

#### Accuracy Metric

#### 2.2.4.4 Discussion

In its presence, the variations caused by these external distortions are amongst the largest of the variations between the faces themselves. In such situations, such as with the Yale B extended Database, the PCA-based approach acquires a noisy mean-face and hence, *Eigenfaces* that captures

## CHAPTER 2. AUTOMATIC FACIAL RECOGNITION

unwanted distortion in its description. Since the images of a particular face, under varying illumination but fixed pose, lie in a 3D linear subspace of the high-dimensional image space (without shadow), *Fisherfaces* work well with the more complex face images, so does PCA.

As expected, overall higher accuracy were obtained on the simpler AT&T dataset. *Eigenfaces* led the 3 algorithms slightly on AT&T with 7 training samples, which was interesting to see. This shows the true power of the PCA-based approach in the ideal case; the lack of unwanted variations in the simple AT&T dataset makes *Eigenfaces* the perfect candidate for it. Eigenfaces score lowest on more complex dataset; PCA does the best, especially with illumination variations; *Fisherfaces* is the runner-up.

For the average time to train, *Eigenfaces* are consistently the slowest; PCA is the quickest overall, as it steadily increases in time with increasing number of training samples; and *Fisherfaces* are consistently quicker than *Eigenfaces*, but always slower than PCA.

Provided an analysis and comparison between a few modern facial recognition techniques. *Fisherfaces* demonstrated top performance in terms of accuracy and performance. PCA achieves a high accuracy at the price of an increased testing time.

### 2.3 Modern-day, data driven deep learning

Hence, the traditional workflow of a machine learning system have the feature extraction and modeling as independent modules (*e.g.*, a system extracts color histograms as the feature and trains a nearest neighbor model on top). In other words, as a schematic or flowchart drawing, the features and the model are separate entities, while deep learning encapsulates the two steps into one—a known benefit of deep learning technology is that the models find the most interesting features at the bottom (*i.e.*, beginning) of the network, that minimize the loss using the model making up the top. The upside here is a built-in mechanism that determines the features of highest interest, opposed to a human having to determine the optimal feature type for a problem, making it such that an expert of a specific domain was typically required for each data type. Largely, it was typically the feature extraction step that demanded specialists with years of knowledge of the quarks and in depth to acquire a means of hand-crafting features. Nowadays, the feature extractor and trained model are one and the same. Of course, domain specific knowledge is still essential in some instances and tasks. However, it is from this that researchers were then enabled to design multimodal systems more frequently; like no other time were there as many SOTA approaches shared across problem domains—attention for text [38], vision [39], speech [40], graph structures [41], and others. The

## CHAPTER 2. AUTOMATIC FACIAL RECOGNITION

same even holds true for major topics such as recurrent neural networks (RNNs) and CNNs, along with fundamental concepts like back-propagation [42] and drop-out [43]. Note that this is not a claim that knowledge has only recently begun to transfer between problem domains in such an explicit manner, as that would be by no means accurate (*e.g.*, bag of words models [44]). However, the lines separating experts that specialize in specific data types seems to be fading away. Furthermore, provided scripts to preprocess the inputs to a neural network (NN), the workings of deep learning technology tends to remain across different data types. Hence, the *black box* encapsulated within a deep model shares tendencies for various signal types.

Modern-day deep models have worked with prominence in automatic face understanding problems. In 2014 Taigman *et al.* of Facebook first proposed using a deep CNN for FR [45]. Over a half of a decade later there has been at least one major contribution in conventional FR each year ever since (*i.e.*, 2015 [46], 2016 [47], 2017 [48], 2018 [49], 2019 [50], and even 2020, the year of facial masks [51]). Details on deep learning advances in FR technology are provided as a part of recent surveys [52, 53].

To briefly summarize, and to a degree needed for a more complete understanding of the work in using facial cues to detect family members, let us better understand the uniqueness to face-based problems. In conventional object recognition, the task to identify a predefined object in imagery (*e.g.*, a cat). For this, one could perceive the problem as boolean (*i.e.*, *is* or *is not* a cat). Alternatively, in a closed-set problem, in a multi-class task, the approach can be founded on determining which class is the instance of (*e.g.*, trained on a set of classes representing *pets*). Regardless, in a supervised setting, the model has access to labeled instances for the object(s) of interest during training. Now, in FR, the multi-class setting for this fine-grained classification problem takes on a different form: given a set of face images for  $C$  classes (*i.e.*, identities), the goal is to train a model to verify a pair of faces for being of the same subject or not. However, the subject in question is not in the *train set*. In fact, the subjects in all *test* photos were never seen by the model prior. By this, FR is inherently a *one shot* learning problem—only one sample per class is assumed, which is the very sample given at inference when asked to determine whether or not it is a match.

Why is it important to understand specifics in settings followed for FR evaluations? How does it being a *one shot* problem change anything? Well, the answers directly relate to overarching goal of a FR system— to learn a mapping that encodes face images as features in such a way that best separates samples of different classes (*i.e.*, identities), while bringing those of the same identity closer together. Hence, a deep model for FR is typically referred to as a face encoder, for the model serves as a feature extractor, and with the encoding as a representation of the face. A pair of faces is

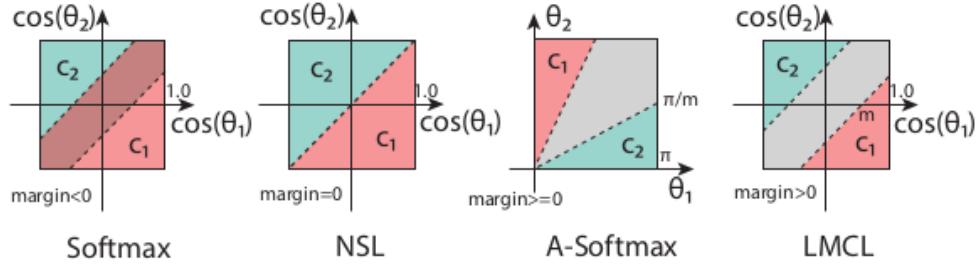


Figure 2.10: **Learned spaces (visualization from [1]).** These schematics are derived from the respective loss function—*left-to-right*: traditional softmax, NSL, ArcFace, and LMCL.

encoded and compared to one another via a distance or similarity metric, where the output of the metric determines whether or not the pair is a true match or not (*i.e.*, does the similarity score surpass a threshold that acts as a decision boundary between *genuine* and *impostor* pairs). This is also the reason many propose Siamese networks [54] as a solution—NNs that share weights for a pair of inputs, which, like in our case, can later be compared in the setting of a verification task. Furthermore, we want a deep face encoder that projects high dimensional images to a highly discriminative feature space (*i.e.*, compactness intra-class and separation inter-class).

### 2.3.1 Loss functions

Like in many machine learning (ML)-based solutions over the past several years is the unquestioned successes with deep learning in automatic FR systems. Inherently, much of the success owes itself to other machine vision solution spaces— network backbones often founded in a more generic problem statement and traditional object were mostly adapted in SOTA FR systems. All the while, it was the loss function that had evolved for face-based model training.

#### 2.3.1.1 Softmax

Softmax is amongst the most popular losses used in deep learning— map the output of the final, topmost layer maps the signal to a score-space by activating the *logits scores* to produce a vector in probability space (*i.e.*, per mathematical axioms that formalize probability theory). Hence, it is here in the deep network that the classification layers are set (*i.e.*, deep learning encapsulates feature extraction and modeling training as a single module, with the features are derived from the input signal at the bottom of the network to then classify from the topmost layers). The output of the layer just below the softmax passes the raw scores, then the vector is mapped as a probability

## CHAPTER 2. AUTOMATIC FACIAL RECOGNITION

vector, which is then compared to the input during training via cross-entropy to produce a loss to back-propagate. Specifically, cross-entropy compares the predicted to the true in classification by comparing the normalized vector to a one-hot encoding representative of the class instance for the sample. Mathematically, given  $N$  training samples of paired data (*i.e.*, image  $x$  and label  $y$ ),

$$L_s = \frac{1}{N} \sum_{i=1}^N -\log(p_i) = \frac{1}{N} \sum_{i=1}^N -\log\left(\frac{e^{f_{y_i}}}{\sum_{j=1}^C e^{f_{y_j}}}\right), \quad (2.7)$$

where the posterior probability  $p$  is conditioned on the predicted class  $c_j$  defined by  $C = \mathbb{R}^k$  for  $|C| = K$ . Then,  $f$ , the activation of the topmost fully-connected layer (*i.e.*, the raw score). Hence, the purpose of the normalized exponential mapping is the softmax operation, while a comparison of bits with that of the true label vector (*i.e.*, one-hot encoding) and the predicted (*i.e.*, normalized between zero and one via softmax) determines the loss that the model parameters are adjusted to minimize.

As emphasized earlier, the underlying goal of encoding faces is to do so such that faces of the same person are close, while faces of different folks are spread far apart. For this, a fair expectation would be an objective that explicitly pushes samples of different classes apart while pulling those of the same closer together—neither is done with the softmax loss. Also, it is susceptible to favoring instances that closely mimics the classes with a majority during training. In other words, imbalanced classes could yield a bias system—a phenomena that is especially of concern in biometrics (more on this in Chapter 8).

Scaled by a weight vector  $W$  and without a bias term (*i.e.*, bias set to zero),

$$\hat{y}_j = W_j^T x = ||W_j|| ||x|| \cos \theta_j. \quad (2.8)$$

Hence, the posterior  $p$  is effected by the angle  $\theta$  (*i.e.*, angle between  $W$  and  $x$ ). For encodings that have a invariant to the norm of  $W$  so the norm is set constant via  $L^1$  and  $L^2$  normalization—during inference, the resulting similarity score of a pair of features only depend on cosine similarity, so with the norm fixed during training the following loss can be determined:

$$L_{ns} = \frac{1}{N} \sum_i -\log\left(\frac{e^{s \cos(\theta_{y_i, i})}}{\sum_j e^{s \cos(\theta_{j, i})}}\right).$$

Now, with the fixed norm  $s$  yields normalized features separable that are separable in angular space—a loss known as the Normalized Version of Softmax Loss (NSL).

# Chapter 3

## Face Detection

### 3.1 Overview

In landmark detection the task is to find the pixel locations in visual media corresponding to points of interest. In face alignment, these points correspond to face parts. For bodies and hands, landmarks correspond to projections of joints on the camera plane [55, 24]. Historically, landmark detection and shape analysis tasks date back decades: from Active Shape Models [56] to Active Appearance Models [57], with the latter proposed to analyze and detect facial landmarks.

Like many other ML-based problems, the task of landmark detection was regained attention with the advancement of deep learning—models capable of encapsulating increasingly tricky views. In other words, the high capacity of deep learning models revamped interest in facial landmark localization— one of the older problems researched in FR [58]. As a result in came a wave of different types of deep neural architectures that pushed SOTA on more challenging datasets. These modern-day networks are trained end-to-end on paired labeled data  $(d, s)$ , where  $d$  is the image and  $s$  are the actual landmark coordinates. Many of these used encoder-decoder style networks to generate feature maps (*i.e.*, heatmaps) to transform into pixel coordinates [59, 60, 61]. The network must be entirely differentiable to train end-to-end. Hence, the layer (or operation) for transforming the  $K$  heatmaps to pixel coordinates must be differentiable [62]. Note that each of the  $K$  heatmaps corresponds to the coordinates of a landmark. Typically, the softargmax operation determines the location of a landmark as the expectation over the generated 2D heatmaps. Thus, metrics like L1 or L2 determine the distance between the actual and predicted coordinates  $\tilde{s}$ , *i.e.*,  $e = \tilde{s} - s$ .

There are two critical shortcomings of the methodology discussed above. (1) These losses only penalize for differences in mean values in coordinate space, and with no explicit penalty for the

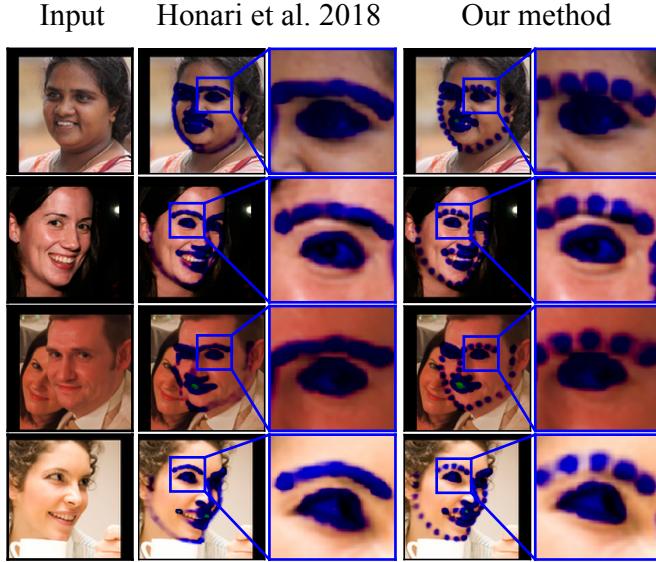


Figure 3.1: **Problem statement.** Heatmaps generated by SAM-based models (middle block) and the proposed LaplaceKL (right block), each with heatmaps on the input images (left) and a zoomed-in view of an eye region (right). These heatmaps are confidence scores (*i.e.*, probabilities) that a pixel is a landmark. Softargmax-based methods generate highly scattered mappings (low certainty), while the same network trained with our loss is concentrated (*i.e.*, high certainty). We further validate the importance of minimizing scatter experimentally (Table 3.2). Best if viewed electronically.

variance of heatmaps. Thus, the generated heatmaps are highly scattered: high variance means low confidence. (2) This family of objectives is entirely dependent on paired training samples (( $\mathbf{d}, \mathbf{s}$ )). However, obtaining high-quality data for this is expensive and challenging. Not only does each sample require several marks, but unintentional, and often unavoidable, labels are of pixel-level marks subject to human error (*i.e.*, inaccurate and imprecise ground-truth labels). All the while, plenty of unlabeled face data are available for free.

## 3.2 Research Contributions

We proposed a practical framework to satisfy the two shortcomings [28]. Specifically, our first contribution alleviates the problem of an unaccounted for spread . For this, we introduce a new loss function that penalizes for the difference in distribution defined by location and scatter (Figure 3.1). Independently, we treat landmarks as random variables with  $\text{Laplace}(\mathbf{s}, 1)$  distributions,

### CHAPTER 3. FACE DETECTION

from which the KL-divergence between the predicted and ground-truth distributions defines the loss. Hence, the goal is to match distributions, parameterized by both a mean and variance, to yield heatmaps of less scatter (*i.e.*, higher confidence). We call this objective the LaplaceKL loss.

Our second contribution in landmark detection was an adversarial training framework. We proposed a method that tackles the problem of paired data requirements by leveraging unlabeled data accessed for free. We treat our landmark detection network as a  $G$  of normalized heatmaps (*i.e.*, probability maps) that pass to the  $D$  to learn to distinguish between the real and fake heatmaps. We could then add large amounts of unlabeled data to further boost the performance of our LaplaceKL-based models. In the end,  $D$  improved the predictive power of the LaplaceKL-based model by injecting unlabeled data into the pipeline during training. Experiments demonstrate that our adversarial training framework complements the proposed LaplaceKL loss (*i.e.*, an increase in unlabeled data results in a decrease in error). We first show the effectiveness of the proposed LaplaceKL loss by claiming SOTA (*i.e.*, first-place) on labeled set and second-to-best (*i.e.*, second-place) without the adversarial training. We then record results for the adversarial training scheme (*i.e.*, leveraging increasing amounts of unlabeled data). Finally, we further improve results using our LaplaceKL loss with more unlabeled data added during training!

Furthermore, we reduced the size of the model by using  $\frac{1}{16}$ ,  $\frac{1}{8}$ ,  $\frac{1}{4}$ , and  $\frac{1}{2}$  the original number of convolution filters, with the smallest costing only 79 Kb on disk. We show an accuracy drop for models trained with the proposed LaplaceKL as far less than the others trained with a softargmax-based loss. So again, it is the case that more unlabeled training data results in less of a performance drop at reduced sizes. It is essential to highlight that variants of our model at or of larger size than 1/8 the original size compare well to the existing SOTA. We claim that the proposed contributions are instrumental for landmark detection models used in real-time production, mobile devices, and other practical purposes.

Our contributions are three-fold: (1) A novel Laplace KL-divergence objective to train landmark localization models that are more certain about predictions; (2) An adversarial training framework that leverages large amounts of unlabeled data during training; (3) Experiments that show our model outperforms recent works in face landmark detection, along with ablation studies that, most notably, reveal our model compares well to SOTA at 1/8 its original size (*i.e.*, <160 Kb) and in real-time (*i.e.*, >20 fps).

### 3.3 Background Information

In this section, we review relevant works on landmark localization and GAN.

#### 3.3.1 Landmark localization

As mentioned, landmark localization (or detection) has been of interest to researchers for decades. At first, most methods were based on Active Shape Models [56] and Active Appearance Models [57]. Then, Cascaded Regression Methods (CRMs) were introduced, which operate sequentially; starting with the average shape, then incrementally shifting the shape closer to the target shape. CRMs offer high speed and accuracy (*i.e.*, >1,000 fps on CPU [63, 64]).

More recently, deep-learning-based approaches have prevailed in the community due to end-to-end learning and improved accuracy. Initial works mimicked the iterative nature of cascaded methods using recurrent convolutional neural networks [60, 65, 66, 67]. Besides, there have been several methods for dense landmark localization [68, 69] and 3D face alignment [70, 71] proposed: all of which are fully-supervised and, thus, require labels for each image.

Nowadays, there is an increasing interest in semi-supervised methods for landmark localization. Recent work used a sequential multitasking method which was capable of injecting labels of two types into the training pipeline, with one type constituting the annotated landmarks and the other type consisting of facial expressions (or hand-gestures) [62]. The authors argued that the latter label type was more easily obtainable, and showed the benefits of using both types of annotations by claiming SOTA on several tasks. Additionally, they explore other semi-supervised techniques (*e.g.*, equivariance loss). In [72], a supervision-by-registration method was proposed, which significantly utilized unlabeled videos for training a landmark detector. The fundamental assumption was that the neighboring frames of the detected landmarks should be consistent with the optical flow computed between the frames. This approach demonstrated a more stable detector for videos, and improved accuracy on public benchmarks.

Landmark localization data resources have significantly evolved as well, with the 68-point mark-up scheme of the MultiPIE dataset [73] widely adopted. Despite the initial excitement for MultiPIE throughout the landmark localization community [74], it is now considered one of the easy datasets captured entirely in a controlled lab setting. A more challenging dataset, Annotated Facial Landmarks in the Wild (AFLW) [75], was then released with up to 21 facial landmarks per face (*i.e.*, occluded or “invisible” landmarks were not marked). Finally, came the 300W dataset made-up of face images from the internet, labeled with the same 68-point mark-up scheme as MultiPIE, and

promoted as a data challenge [76]. Currently, 300W is among the most widely used benchmarks for facial landmark localization. In addition to 2D datasets, the community created several datasets annotated with 3D keypoints [77].

### 3.3.2 GANs

were recently introduced [78], quickly becoming popular in research and practice. GANs have been used to generate images [79] and videos [80, 81], and to do image manipulation [82], text-to-image[83], image-to-image [84], video-to-video [85] translation and re-targeting [86].

An exciting feature of GANs is the ability to transfer visual media across different domains. Thus, various semi-supervised and domain-adaptation tasks adopted GANs [87, 88, 89, 90]. Many have leveraged synthetic data to improve model performance on real data. For example, a GAN transferred images of human eyes from the real domain to bootstrap training data [89]. Other researchers used them to synthetically generate photo-realistic images of outdoor scenes, which also aided in bettering performance in image segmentation [88]. Sometimes, labeling images captured in a controlled setting is manageable (*i.e.*, versus an uncontrolled setting). For instance, 2D body pose annotations were available *in-the-wild*, while 3D annotations mostly were for images captured in a lab setting. Therefore, images with 3D annotations were used in adversarial training to predict 3D human body poses as seen *in-the-wild* [90]. [87] formulated one-shot recognition as a data imbalance problem and augmented additional samples in the form of synthetic embeddings.

Our work differs from these others in several ways. Firstly, a majority, if not all, used a training objective that only accounts for the location of landmarks [62, 65, 67], *i.e.*, no consideration for variance (*i.e.*, confidence). Thus, landmarks distributions have been assumed to be describable with a single parameter (*i.e.*, a mean). Networks trained this way yield an uncertainty about the prediction, while still providing a reasonable location estimate. To mitigate this, we explicitly parametrize the distribution of landmarks using location and scale. For this, we propose a KL-divergence based loss to train the network end-to-end. Secondly, previous works used GANs for domain adaptation in some fashion. In this work, we do not perform any adaptation between domains as in [88, 89], nor do we use any additional training labels as in [62]. Specifically, we have  $D$  do the quality assessment on the predicted heatmaps for a given image. The resulting gradients are used to improve the ability of the generator to detect landmarks. We show that both contributions improve accuracy when used separately. Then, the two contributions together boost SOTA results.

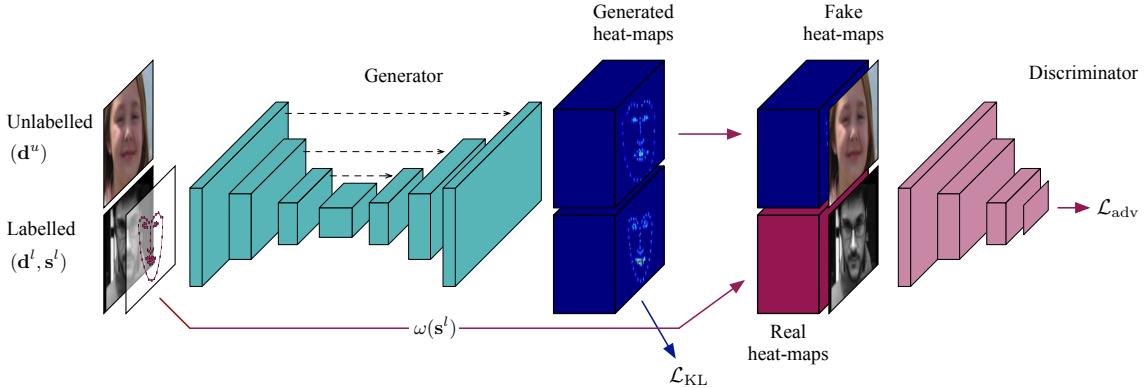


Figure 3.2: **Our semi-supervised framework for landmark detection.** The labeled and unlabeled branches are marked with blue and red arrows, respectfully. Given an input image,  $G$  produces  $K$  heatmaps, one for each landmark. Labels are used to generate real heatmaps as  $\omega(\mathbf{s}^l)$ .  $G$  produces fake samples from unlabeled data. Source images are concatenated on heatmaps and passed to  $D$ .

### 3.4 Laplace Landmark Localizer

Our training framework utilizes both labeled and unlabeled data during training. Shown in Figure 3.2 are the high-level graphical depiction of cases where labels are available (blue arrows) and unavailable (red arrows). Notice the framework has two branches, supervised (Eq. 3.3) and unsupervised (Eq. 3.7), where only the supervised (blue arrow) uses labels to train. Next, are details for both branches.

#### 3.4.1 Fully Supervised Branch

We define the joint distribution of the image  $\mathbf{d} \in \mathbb{R}^{h \times w \times 3}$  and landmarks  $\mathbf{s} \in \mathbb{R}^{K \times 2}$  as  $p(\mathbf{d}, \mathbf{s})$ , where  $K$  is the total number of landmarks. The form of the distribution  $p(\mathbf{d}, \mathbf{s})$  is unknown; however, joint samples are available when labels are present (*i.e.*,  $(\mathbf{d}, \mathbf{s}) \sim p(\mathbf{d}, \mathbf{s})$ ). During training, we aim to learn a conditional distribution  $q_\theta(\mathbf{s}|\mathbf{d})$  modeled by a neural network with parameters  $\theta$ . Landmarks are then detected done by sampling  $\tilde{\mathbf{s}} \sim q_\theta(\mathbf{s}|\mathbf{d})$ . We now omit parameters  $\theta$  from notation for cleaner expressions. The parameter values are learned by maximizing the likelihood that the process described by the model did indeed produce the data that was observed, *i.e.*, trained by minimizing the following loss function w.r.t. its parameters:

$$\mathcal{L}(\theta) = \mathbb{E}_{(\mathbf{d}, \mathbf{s}) \sim p(\mathbf{d}, \mathbf{s})} \|\tilde{\mathbf{s}} - \mathbf{s}\|_2. \quad (3.1)$$

### CHAPTER 3. FACE DETECTION

Alternatively, it is possible to train a neural network to predict normalized probability maps(*i.e.*, heatmaps):  $\tilde{\mathbf{h}} \sim q(\mathbf{h}|\mathbf{d})$ , where  $\mathbf{h} \in \mathbb{R}^{K \times h \times w}$  and each  $\mathbf{h}_k \in \mathbb{R}^{h \times w}$  represents a normalized probability map for landmark  $k$ , where  $k = 1 \dots K$ . To get the pixel locations, one could perform the argmax operation over the heatmaps by setting  $\tilde{\mathbf{s}} = \text{argmax}(\tilde{\mathbf{h}})$ ). However, this operation is not differentiable and, therefore, unable to be trained end-to-end.

A differentiable variant of argmax (*i.e.*, softargmax [91]) was recently used to localize landmarks [62]. For the 1D case, the softargmax operation is expressed

$$\begin{aligned}\text{softargmax}(\beta\mathbf{h}) &= \sum_x \text{softmax}(\beta\mathbf{h}_x) \cdot x \\ &= \sum_x \frac{e^{\beta\mathbf{h}_x}}{\sum_j e^{\beta\mathbf{h}_j}} \cdot x \\ &= \sum_x p(x) \cdot x = \mathbb{E}_{\mathbf{h}}[x],\end{aligned}\tag{3.2}$$

where  $\mathbf{h}_x$  is the predicted probability mass at location  $x$ ,  $\sum_j e^{\beta\mathbf{h}_j}$  is the normalization factor, and  $\beta$  is the temperature factor controlling the predicted distribution [91]. Coordinates are in boldface (*i.e.*,  $\mathbf{x} = (x_1, x_2)$ ), and write 2D softargmax operation as  $\tilde{\mathbf{s}} = \mathbb{E}_{\mathbf{h}}[\mathbf{x}]$  with  $\mathcal{L}_{\text{SAM}} = \mathcal{L}(\theta)$ .

Essentially, the softargmax operation is the expectation of the pixel coordinate over the selected dimension. Hence, the softargmax-based loss assumes the underlying distribution is describable by just its mean (*i.e.*, location), regardless of how sure a prediction, the objective then is to match mean values. To avoid cases in which the trained model is uncertain about the predicted mean, while still yielding a low error, we parameterize the distribution using  $\{\mu, \sigma\}$ , where  $\mu$  is the mean or the location and  $\sigma$  is the variance or the scale, respectfully, for the selected distribution.

We want the model to be certain about the predictions (*i.e.*, a small variance or scale). We consider two parametric distributions  $\text{Gaussian}(\mu, \sigma)$  and  $\text{Laplace}(\mu, b)$  with  $\sigma^2 = \mathbb{E}_{\mathbf{h}}[(\mathbf{x} - \mathbb{E}_{\mathbf{h}}[\mathbf{x}])^2]$  and  $b = \mathbb{E}_{\mathbf{h}}[|\mathbf{x} - \mathbb{E}_{\mathbf{h}}[\mathbf{x}]|]$ . We define a function  $\tau(\tilde{\mathbf{h}})$  to compute the scale (or variance) of the predicted heatmaps  $\tilde{\mathbf{h}}$  using the location, where the locations are now the expectation of being a landmark in the heatmap space. Thus,  $\tau(\tilde{\mathbf{h}}) = \sum p(\mathbf{x}) \|\mathbf{x} - \tilde{\mathbf{s}}\|_{\alpha}^{\alpha}$ , where  $\tilde{\mathbf{s}} = \mathbb{E}_{\mathbf{h}}[\mathbf{x}]$ ,  $\alpha = 1$  for Laplacian, and  $\alpha = 2$  for Gaussian. Thus,  $\tilde{\mathbf{s}}$  and  $\tau(\tilde{\mathbf{h}})$  are used to parameterize a Laplace (or Gaussian) distribution for the predicted landmarks  $q(\mathbf{h}|\mathbf{d})$ .

With the true conditional distribution of the landmarks as  $p(\mathbf{s}|\mathbf{d})$ , the objective is

$$\mathcal{L}_{\text{KL}} = \mathbb{E}_{(\mathbf{d}, \mathbf{s}) \sim p(\mathbf{d}, \mathbf{s})} \left[ D_{\text{KL}}(q(\mathbf{s}|\mathbf{d}) || p(\mathbf{s}|\mathbf{d})) \right],\tag{3.3}$$

---

**Algorithm 1:** Training the proposed model.
 

---

```

Data:  $\{(\mathbf{d}_i^l, \mathbf{s}_i^l)\}_{i=1,\dots,n}, \{(\mathbf{d}_i^u)\}_{i=1,\dots,m}$ 
 $\theta_D, \theta_G \leftarrow$  initialize network parameters
while  $t \leq T$  do
     $(\mathbf{D}_t^l, \mathbf{S}_t^l) \leftarrow$  sample mini-batch from labeled data
     $(\mathbf{D}_t^u) \leftarrow$  sample mini-batch from unlabeled data
     $\mathbf{H}_{\text{fake}} \leftarrow G(\mathbf{D}_t^u)$ 
     $\mathbf{H}_{\text{real}} \leftarrow \omega(\mathbf{S}_t^l)$ 
     $\mathcal{L}_{\text{adv}} \leftarrow \log D([\mathbf{D}_t^l, \mathbf{H}_{\text{real}}]) + \log(1 - D([\mathbf{D}_t^u, \mathbf{H}_{\text{fake}}]))$ 
     $\mathcal{L}_G \leftarrow$  compute loss using Eq. 3.2 or Eq. 3.3
    // update model parameters
     $\theta_D \xleftarrow{+} -\nabla_{\theta_D} \mathcal{L}_{\text{adv}}$ 
     $\theta_G \xleftarrow{+} -\nabla_{\theta_G} (\mathcal{L}_G - \lambda \mathcal{L}_{\text{adv}})$ 
end
    
```

---

where  $D_{\text{KL}}$  is the KL-divergence. We assumed a true distribution for the case of Gaussian (*i.e.*,  $\text{Gaussian}(\mu, 1)$ , where  $\mu$  is the ground-truth locations of the landmarks). For the case with Laplace, we sought  $\text{Laplace}(\mu, 1)$ . KL-divergence conveniently has a closed-form solution for this family of exponential distributions [92]. Alternatively, sampling yields an approximation. The blue arrow in Figure 3.2 represent the labeled branch of the framework.

Statistically speaking, given two estimators with different variances, we would prefer one that has a smaller variance (see [93] for an analysis of the bias-variance trade-off). A lower variance implies higher confidence in the prediction. To this end, we found an objective measuring distance between distributions is accurate and robust. The neural network must satisfy an extra constraint on variance and, thus, yields predictions of higher certainty. See higher confident heatmaps in Figure 3.1 and Figure 3.3. The experimental evaluation further validates this (Table 3.2 and Table 3.3). Also, Figure 3.4 shows sample results.

### 3.4.2 Unsupervised Branch

The previous section discusses several objectives to train the neural network with the available paired or fully labeled data (*i.e.*,  $(\mathbf{d}^l, \mathbf{s}^l)$ ). We denote data samples with the superscript  $l$  to distinguish them from unpaired or unlabeled data  $(\mathbf{d}^u)$ . In general, it is difficult for a human

to label many images with landmarks. Hence, unlabeled data are abundant and easier to obtain, which calls for capitalizing on this abundant data to improve training. In order to do so, we adapt the adversarial learning framework for landmark localization. We treat our landmarks predicting network as a generator ( $G$ ),  $G = q(\mathbf{h}|\mathbf{d})$ ; discriminator ( $D$ ) takes the form  $D([\mathbf{d}, \mathbf{h}])$ , where  $[\cdot, \cdot]$  is a tensor concatenation operation. We define the real samples for  $D$  as  $\{\mathbf{d}^l, \mathbf{h} = \omega(\mathbf{s}^l)\}$ , where  $\omega(\cdot)$  generates the true heatmaps given the locations of the ground-truth landmarks. Fake samples are given by  $\{\mathbf{d}^u, \tilde{\mathbf{h}} \sim q(\mathbf{h}|\mathbf{d}^u)\}$ . We then define the min-max objective for landmark detection as:

$$\min_G \max_D \mathcal{L}_{\text{adv}}(D, G), \quad (3.4)$$

where  $\mathcal{L}_{\text{adv}}(D, G)$  writes as:

$$\begin{aligned} & \mathbb{E}_{(\mathbf{d}^l, \mathbf{s}^l) \sim p(\mathbf{d}, \mathbf{s})} \left[ \log D([\mathbf{d}^l, \omega(\mathbf{s}^l)]) \right] + \\ & \mathbb{E}_{(\mathbf{d}^u) \sim p(\mathbf{d})} \left[ \log(1 - D([\mathbf{d}^u, G(\mathbf{d}^u)])) \right]. \end{aligned} \quad (3.5)$$

With this, provided an input image, the goal of  $D$  is to learn to decipher between the real and fake heatmaps from appearance. The goal of  $G$  is to produce fake heatmaps that closely resemble the real. Within this framework,  $D$  intends to provide additional guidance for  $G$  by learning from labeled and unlabeled data. The objective, Eq. 3.4, is solved using alternating updates.

### 3.4.3 Training

We fused the softargmax-based and adversarial losses as

$$\min_G \left( \max_D (\lambda \cdot \mathcal{L}_{\text{adv}}(G, D)) + \mathcal{L}_{\text{SAM}}(G) \right), \quad (3.6)$$

with the KL-divergence version of the objective defined as:

$$\min_G \left( \max_D (\lambda \cdot \mathcal{L}_{\text{adv}}(G, D)) + \mathcal{L}_{\text{KL}}(G) \right), \quad (3.7)$$

with the weight for the adversarial loss  $\lambda = 0.001$ . This training objective includes both labeled and unlabeled data in the formulation. In the experiments, we show that this combination significantly improves the accuracy of our approach. We also argue that the softargmax-based version cannot fully utilize the unlabeled data since the predicted heatmaps differ too much from the *real* heatmaps. See Algorithm 1 for the training procedure for  $T$  steps of the proposed model. We show the unlabeled branch of the framework is shown graphically in red arrows (Figure 3.2).

Table 3.1: **Architecture of generator ( $G$ ).** Layers with size and number of filters (*i.e.*,  $h \times w \times n$ ). DROP, MAX, and UP are dropout (probability 0.2), max-pooling (stride 2), and bilinear upsampling ( $2x$ ), respectively. Note, skip connections about the bottleneck: coarse-to-fine, connecting encoder ( $E_{ID}$ ) to decoder ( $D_{ID}$ ) by concatenating feature channels and fusing. Number of feature preserved at all but top two layers (*i.e.*, transform → features →  $K$  heatmaps). Padded to match sizes listed.

	Layers	Tensor Size
Input	RGB image, no data augmentation	$80 \times 80 \times 3$
Conv( $E_1$ )	$3 \times 3 \times 64$ , LReLU, DROP, MAX	$40 \times 40 \times 64$
Conv( $E_2$ )	$3 \times 3 \times 64$ , LReLU, DROP, MAX	$20 \times 20 \times 64$
Conv( $E_3$ )	$3 \times 3 \times 64$ , LReLU, DROP, MAX	$10 \times 10 \times 64$
Conv( $E_4$ )	$3 \times 3 \times 64$ , LReLU, DROP, MAX	$5 \times 5 \times 64$
Conv( $D_4$ )	$1 \times 1 \times 64 + E_4$ , LReLU, DROP, UP	$10 \times 10 \times 128$
Conv( $D_F$ )	$5 \times 5 \times 128$ , LReLU	$20 \times 20 \times 128$
Conv( $D_3$ )	$1 \times 1 \times 64 + E_3$ , LReLU, DROP, UP	$20 \times 20 \times 128$
Conv( $D_F$ )	$5 \times 5 \times 128$ , LReLU, DROP	$40 \times 40 \times 128$
Conv( $D_2$ )	$1 \times 1 \times 64 + E_2$ , LReLU, DROP, UP	$40 \times 40 \times 128$
Conv( $D_F$ )	$5 \times 5 \times 128$ , LReLU, DROP	$80 \times 80 \times 128$
Conv( $D_1$ )	$1 \times 1 \times 64 + E_1$ , LReLU, DROP, UP	$80 \times 80 \times 128$
Conv( $D_F$ )	$5 \times 5 \times 128$ , LReLU, DROP	$80 \times 80 \times 128$
Conv( $D_F$ )	$1 \times 1 \times 68$ , LReLU, DROP	$80 \times 80 \times 68$
Output	$1 \times 1 \times 68$	$80 \times 80 \times 68$

### 3.4.4 Implementation

We follow the ReCombinator network (RCN) initially proposed in [94]. Specifically, we use a 4-branch RCN as our base model, with input images and output heatmaps of size  $80 \times 80$ . Convolutional layers of the encoder consist of 64 channels, while the convolutional layers of the decoder output 64 channels out of the 128 channels at its input (*i.e.*, 64 channels from the previous layer concatenated with the 64 channels skipped over the bottleneck via branching). We applied Leaky-ReLU, with a negative slope of 0.2, on all but the last convolution layer. See Table 3.1 for details on the generator architecture. Drop-out followed this, after all but the first and last activation. We use Adam optimizer with a learning rate of 0.001 and weight decay of  $10^{-5}$ . In all cases, networks were trained from scratch, with no data augmentation or ‘training tricks’.

$D$  was a 4-layered PatchGAN [95]. Before each convolution layer Gaussian noise ( $\sigma = 0.2$ ) was added [81], and then batch-normalization (all but the top and bottom layers) and Leaky-ReLU with a negative slope of 0.2 (all but the top layer). The original RGB image was stacked on top of the  $K$  heatmaps from  $G$  and fed as the input of  $D$  (Figure 3.2). Thus,  $D$  takes in  $(K + 3)$  channels. We set  $\beta = 1$  for 3.2. Pytorch was used to implement the entire framework. An important note to make is that models optimized with Laplace distribution consistently outperformed the Gaussian-based.

Table 3.2: **Quantitative results.** NMSE on AFLW and 300W normalized by the square root of BB area and interocular distance, respectfully.

	AFLW	300W		
		Common	Challenge	Full
SDM [96]	5.43	5.57	15.40	7.52
LBF [63]	4.25	4.95	11.98	6.32
MDM [65]	-	4.83	10.14	5.88
TCDCN [97]	-	4.80	8.60	5.54
CFSS [98]	3.92	4.73	9.98	5.76
CFSS [99]	2.17	4.36	7.56	4.99
RCSR [67]	-	4.01	8.58	4.90
RCN+ (L+ELT) [62]	<b>1.59</b>	4.20	7.78	4.90
CPM + SBR [72]	2.14	3.28	7.58	4.10
<hr/>				
Softargmax	2.26	3.48	7.39	4.25
Softargmax+D(10K)	-	3.34	7.90	4.23
Softargmax+D(30K)	-	3.41	7.99	4.31
Softargmax+D(50K)	-	3.41	8.06	4.32
Softargmax+D(70K)	-	3.34	8.17	4.29
<hr/>				
LaplaceKL [28]	1.97	3.28	7.01	4.01
LaplaceKL+D(10K)	-	3.26	6.96	3.99
LaplaceKL+D(30K)	-	3.29	6.74	3.96
LaplaceKL+D(50K)	-	3.26	<b>6.71</b>	3.94
LaplaceKL+D(70K)	-	<b>3.19</b>	6.87	<b>3.91</b>

For instance, our LaplaceKL baseline had a NMSE of 4.01 on 300W, while Gaussian-based got 4.71. Thus, the sharper, “peakier” Laplace distribution proved to be more numerically stable under current network configuration, as Gaussian required a learning rate of a magnitude smaller to avoid vanishing gradients. Indeed, we used Laplace.

## 3.5 Experiments

We evaluated the proposed on two widely used benchmark datasets for face alignment. No data augmentation techniques used when training our models nor was the learning rate dropped: this leaves no ambiguity into whether or not the improved performance came from training tricks or the learning component itself. All results for the proposed were from models trained for 200 epochs.

We next discuss the metric used to evaluate performance, NMSE, with differences between datasets in the normalization factor. Then, the experimental settings, results, and analysis for each dataset are covered separately. Finally, ablation studies show characterizations of critical hyperparameters and, furthermore, the robustness of the proposed LaplaceKL+D(70K) with a comparable performance with just 1/8 the number of feature channels and >20 fps.

### 3.5.1 Metric

Per convention [76, 77, 100], NMSE, a normalized average of euclidean distances, was used. Mathematically speaking:

$$\text{NMSE} = \sum_{k=1}^K \frac{\|s_k - \tilde{s}_k\|_2}{K \times d}, \quad (3.8)$$

where the number of visible landmarks set as  $K$ ,  $k = \{1, 2, \dots, K\}$  are the indices of the visible landmark, the normalization factor  $d$  depends on the face size, and  $s_k \in \mathbb{R}^2$  and  $\tilde{s}_k \in \mathbb{R}^2$  are the ground-truth and predicted coordinates, respectfully. The face size  $d$  ensured that the NMSE scores across faces of different size were fairly weighted. Following predecessors, NMSE was used to evaluate both datasets, except with different points referenced to calculate  $d$ . The following subsections provide details for finding  $d$ .

### 3.5.2 300W + MegaFace

The 300W dataset is amongst the most popular datasets for face alignment. It has 68 visible landmarks (*i.e.*,  $K = 68$ ) for 3,837 images (*i.e.*, 3,148 training and 689 test). We followed the protocol of the 300W challenge [76] and evaluated using NMSE (Eq. 3.8), where  $d$  is set as the interocular distance (*i.e.*, distance between outer corners of the eyes). Per convention, we evaluated different subsets of 300W (*i.e.*, *common* and *challenge*, which together form *full*).

We compared the performance of the proposed objective trained in a semi-supervised fashion. During training, 300W dataset made-up the labeled data (*i.e.*, *real*), and a random selection from MegaFace provided the unlabeled data (*i.e.*, *fake*) [101]. MTCNN<sup>1</sup> was used to detect five landmarks (*i.e.*, eye pupils, corners of the mouth, and middle of nose and chin) [102], which allowed for similar face crops from either dataset. Specifically, we extended the square hull that enclosed the five landmarks by  $2 \times$  the radii in each direction. In other words, the smallest bounding box spanning the 5 points (*i.e.*, the outermost points lied on the parameter), and then transformed from rectangles-to-squares with sides of length  $2 \times \max(\text{height}, \text{width})$ . Note that the midpoint of the original rectangle was held constant to avoid shift translations (*i.e.*, rounded up a pixel if the radius was even and extended in all directions).

The LaplaceKL+D(70K) model obtained SOTA on 300W, yielding the lowest error on 300W (Table 3.2 (300W columns)). LaplaceKL and softargmax with  $N$  unlabeled faces, s+D( $N$ ).

---

<sup>1</sup><https://github.com/davidsandberg/facenet>

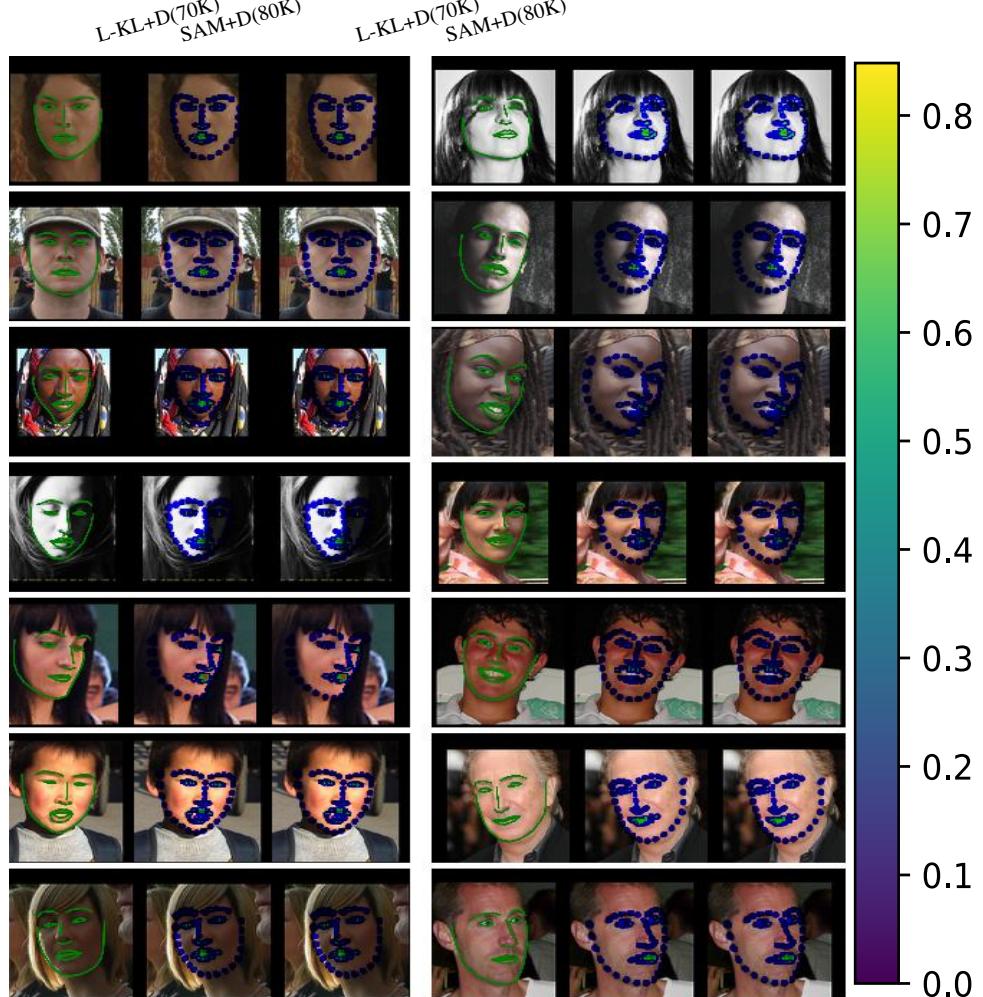


Figure 3.3: **Random samples (300W).** Heatmaps predicted by our LaplaceKL+D(70K) (middle, *i.e.*, L-KL+D(70K)) and softargmax+D(70K) (right, *i.e.*, SAM+D(70K)) alongside face images with ground-truth sketched on the face (left). For this, colors were set by value for the  $K$  heatmaps generated for each landmark (*i.e.*, range of [0, 1] as shown in color bar), and then were superimposed on the original face. Note that the KL-divergence loss yields predictions of much greater confidence and, hence, produced separated landmarks when visualized heatmap space. In other words, the proposed has minimal spread about the mean, as opposed to the softargmax-based model with heatmaps with individual landmarks smudged together. Best viewed electronically.



Figure 3.4: **Qualitative results.** Random samples of landmarks predicted using LaplaceKL (white), with the ground truth drawn as line segments (red). Notice the predicted points tend to overlap with the ground-truth. Best viewed in color. Zoom-in for greater detail.

denote the models trained with unlabeled data, where  $N$  representing the number of unlabeled images added from MegaFace.

First, notice that LaplaceKL trained without unlabeled data still achieved SOTA. The LaplaceKL-based models then showed relative improvements with more unlabeled data added. The softargmax-based models cannot fully take advantage of the unlabeled data without minimizing for variance (*i.e.*, generates heatmaps of less confidence and, thus, more spread). Our LaplaceKL, on the other hand, penalizes for spread (*i.e.*, scale), making the job of  $D$  more challenging. As such, LaplaceKL-based models benefit from increasing amounts of unlabeled data.

Also, notice the largest gap between the baseline models [72] and our LaplaceKL+D(70K) model on the different sets of 300W. Adding more unlabeled helps more (*i.e.*, LaplaceKL versus LaplaceKL+D(70K) improvement is about 2.53%). However, it is essential to use samples not covered in the labeled set. To demonstrate this, we set the *real* and *fake* sets to 300W (*i.e.*,  $\mathbf{d}^l = \mathbf{d}^u$  in the second term of Eq. 3.7). NMSE results for this experiment are listed as follows: LaplaceKL+D(300W) 4.06 (baseline– 4.01) and softargmax+D(300W) 4.26 (baseline– 4.24). As hypothesized, all the information from the labeled set had already been extracted in the supervised branch, leaving no benefit of using the same set in the unsupervised branch. Therefore, more unlabeled data yields more hard negatives to train with, which improves the accuracy of the rarely seen samples (Table 3.2 (300W *challenge* set)). Our best model was  $\approx 2.7\%$  better than [72] on easier

Table 3.3: **Ablation.** NMSE on 300W (full set) for networks trained with fewer channels in each convolutional layer by 1/16, 1/8, 1/4, 1/2, and unmodified in size (*i.e.*, the original) listed from left-to-right. We measured performance with a 2.8GHz Intel Core i7 CPU.

	Number of parameters, millions				
	0.0174	0.0389	0.1281	0.4781	1.8724
Softargmax	9.79	6.86	4.83	4.35	4.25
Softargmax+D(70K)	9.02	6.84	4.85	4.38	4.29
LaplaceKL	7.38	5.09	4.39	4.04	4.01
LaplaceKL+D(70K)	<b>7.01</b>	<b>4.85</b>	<b>4.30</b>	<b>3.98</b>	<b>3.91</b>
Storage (MB)	0.076	0.162	0.507	1.919	7.496
Speed (fps)	26.51	21.38	16.77	11.92	4.92

samples (*i.e.*, *common*),  $\approx 4.7\%$  better on average (*i.e.*, *full*), and, moreover,  $\approx 9.8\%$  better on the more difficult (*i.e.*, *challenge*),  $\approx 4.7\%$  better on average (*full*), and, moreover,  $\approx 9.8\%$  better on the more difficult (*challenge*). These results further highlight the advantages of training with the proposed LaplaceKL loss, along with the adversarial training framework.

Additionally, the adversarial framework boosted our 300W baseline (*i.e.*, more unlabeled data yields a lower NMSE). Specifically, we demonstrated this by pushing SOTA of the proposed on 300W from a NMSE of 4.01 to 3.91 (*i.e.*, no unlabeled data to 70K unlabeled pairs, respectfully). There were boosts at each step size of *full* (*i.e.*, larger  $N \rightarrow$  NMSE).

We randomly sampled the unlabeled for LaplaceKL+D(70K) and softargmax+D(70K) to visualize predicted heatmaps (Figure 3.3). In each case, the heatmaps produced by the softargmax-based models spread wider, explaining the worsened quantitative scores (Table 3.2). Our loss and adversarial learning scheme yield higher probable pixel location (*i.e.*, a more concentrated predicted heatmaps). For most images, the heatmaps generated by models trained with the LaplaceKL loss have distributions for landmarks of more confidence and properly distributed: LaplaceKL+D(70K) yielded heatmaps that vary  $\pm 1.02$  pixels from the mean, while softargmax+D(70K) has a variation of  $\pm 2.59$ . Learning the landmark distributions with our LaplaceKL loss is conceptually and theoretically intuitive (Figure 3.1), and experimentally proven (Table 3.2).

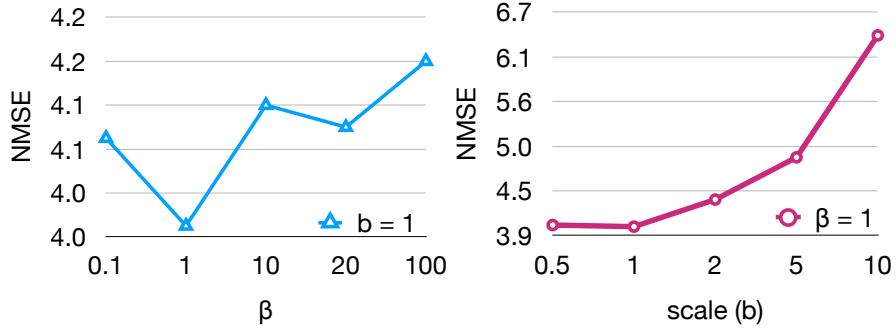


Figure 3.5: **Ablation.** Results of ablation study on LaplaceKL.

### 3.5.3 The Annotated Facial Landmarks in the Wild (AFLW) dataset

We evaluated the LaplaceKL loss on the AFLW dataset [75]. AFLW contains 24,386 faces with up to 21 landmarks annotations and 3D head pose labels. Following [62], 20,000 faces were used for training with the other 4,386 for testing. We ignored the two landmarks for the left and right earlobes, leaving up to 19 landmarks per face [72].

Since faces of AFLW have such variety head poses, most faces have landmarks out of view (*i.e.*, missing). Thus, most samples were not annotated with the complete 19 landmarks, meaning that it does not allow for a constant sized tensor (*i.e.*, *real* heatmaps) for the adversarial training. Therefore, we compared the softargmax and KL-based objectives with existing SOTA. The face size  $d$  for the NMSE was the square root of the bounding box hull [77].

Our LaplaceKL-based model was comparable to SOTA (*i.e.*, RCN+ (L+ELT) [62]) on the larger, more challenging AFLW dataset while outperforming all others. It is essential to highlight here that [62] puts great emphasis on data augmentation, while we do not apply any. Also, since landmarks are missing in some samples (*i.e.*, no common reference points exist across all samples), we were unable to prepare faces for our semi-supervised component– a subject for future work.

### 3.5.4 Ablation Study

The error is next measured as a function of model size (Table 3.3), along with different  $\beta$  values (Eq. 3.2) and scales  $b$  used to parameterize the Laplacian (Figure 3.5). The latter characterizes the baseline and supports the values used for these hyper-parameters, while the former reveals a critical characteristic for the practicality of the proposed.

Specifically, we decreased the model size by reducing the number of channels at each

convolutional layer by factors of 2. The softargmax-based model worsened by about 13% and 59% in NMSE at a  $\frac{1}{4}$  and  $\frac{1}{8}$  the channel count, respectfully (*i.e.*, 4.25 → 6.86 and 9.79). LaplaceKL, on the other hand, decreased by about 24% with an 8<sup>th</sup> and 59% with a 16<sup>th</sup> the number of channels (*i.e.*, 4.01 → 5.09 and 7.38, respectfully). Our model trained with unlabeled data (*i.e.*, LaplaceKL+D(70K)) dropped just about 21% and 57% at factors of 8 and 16, respectfully (*i.e.*, 3.91 → 4.85 and 7.01). LaplaceKL+D(70K) proved best with reduced sizes: with <0.040M parameters, it still compares to previous SOTA [62, 67, 99], which is a clear advantage. For instance, SDM [96], requires 1.693M parameters (25.17MB)<sup>2</sup> for 7.52 in NMSE (300W *full*); our smallest and next-to-smallest got 7.01 and 4.85 with 0.174M (0.076 MB) and 0.340M (0.166 MB) weights.

The model also speeds up with fewer channels (*i.e.*, to train and at inference). For instance, the model reduced by a factor of 16 processes 26.51 frames per second (fps) on a CPU of Macbook Pro (*i.e.*, 2.8GHz Intel Core i7), with the original running at 4.92 fps. Our best LaplaceKL-based model proved robust to size reduction: 4.85 NMSE at 21.38 fps when reduced by 1/8.

### 3.6 Discussion

We demonstrated the benefits of the proposed LaplaceKL loss and leveraging unlabeled data in an adversarial training framework. Hypothetically and empirically, we showed the importance of penalizing a landmark predictor’s uncertainty. Thus, training with the proposed objective yields predictions of higher confidence, outperforming previous SOTA methods. We also revealed the benefits of adding unlabeled training data to boost performance via adversarial training. In the end, our model performs SOTA on all three splits of the renown 300W (*i.e.*, *common*, *challenge*, and *full*), and second-to-best on the AFLW benchmark. Also, we demonstrate the robustness of the proposed by significantly reducing the number of parameters. Specifically, with 1/8 the number of channels (*i.e.*, <170Kb on disk), the proposed still yields an accuracy comparable to the previous SOTA in real-time (*i.e.*, 21.38 fps). Thus, the contributions of the proposed framework are instrumental for models intended for use in real-world production.

---

<sup>2</sup>[https://github.com/tntrung/sdm\\_face\\_alignment](https://github.com/tntrung/sdm_face_alignment)

## **Part II**

# **Visual Kinship Recognition of Families In the Wild**

## Chapter 4

# Visual Kinship Recognition

### 4.1 Overview

About a decade ago, pioneers in visual kinship recognition research published the seminar work in detecting family relationships with face images [34]. Before discussing our specific research contributions in kinship recognition technology (Chapter 5), we first look back at the last decade of research done by the community as a whole: review the key milestones that brought us to the state in technology for where we are today. Furthermore, let us highlight the key challenges, practical use-cases, and promising future directions for research. By doing so, we will have paved the way to diving into the details of the many related efforts:

- Our purpose and motivations;
- The technical novelty and the ways at which it fits into the bigger picture;
- Our hopes and beliefs for which the work we had done as part of this dissertation could and should be considered by other researchers;
- Whether a junior scholar looking for a problem to hone in on as part of a dissertation;
- Experts that have published in the automatic kinship recognition problem space.

In any case, there are great benefits to reap from the advancement of kinship recognition—it has a multitude of practical and scholarly uses. Relationships provide rich information in sociology, anthropology, and genetics; privacy protections and concerns, along with potential use-cases that can be found in social media, personal discovery, entertainment, and more. Besides its entrepreneurial

## CHAPTER 4. VISUAL KINSHIP RECOGNITION

value, visual kinship recognition has significant non-commercial (or humane) value as well. For instance, in cases of missing children, reconnecting families split across refugee camps, border control and customs, criminal investigations, ancestral-based studies, and even genome-based research. Socially, family gives a sense of belonging (*i.e.*, membership, connection). Per Furstenberg,

... important function of family systems receives far less attention in the literature than it merits: The family ... social arrangement responsible for giving its members a sense of identity and shared belonging ... not only those inside the natal family household but also among relations living elsewhere as well [103].

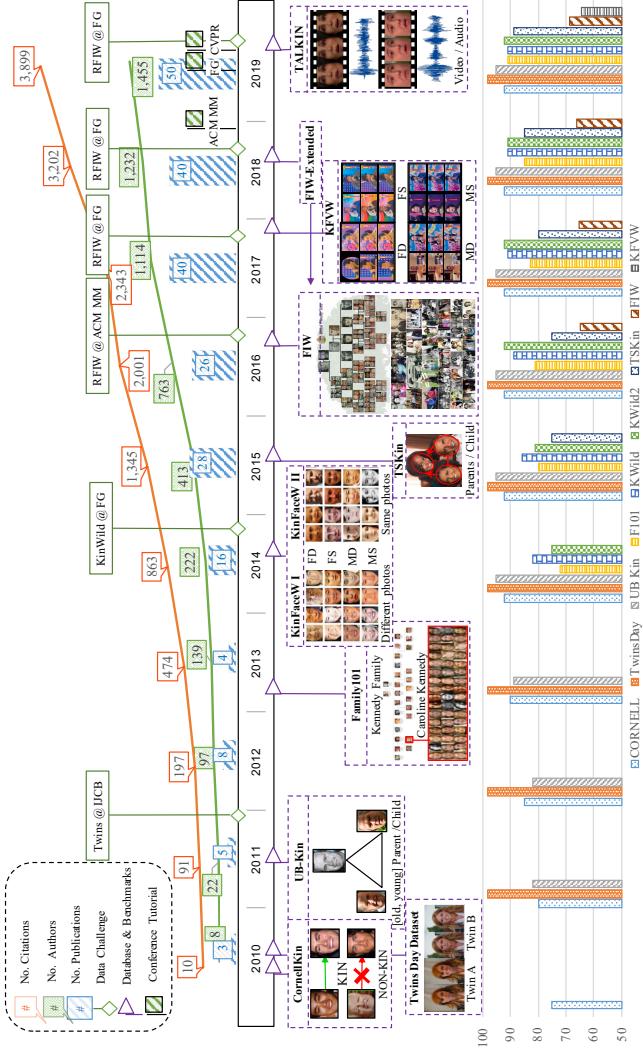
Hence, a recent surge in many seeking out their pedigree. With an abundance of visual data online, familial resources can benefit.

We now review the state of automatic kinship recognition after the first decade of research - with emphasis on the milestones that led us up to now (Fig. Figure 4.2). Furthermore, we reflect on the problem statements of the different tasks to establish clear definitions and an understanding across the domain in a consistent manner. For this, we aim to use consistent terminology, assess the practical usefulness (or lack thereof), and highlight any outstanding challenges and obstacles that prevent the transition of visual kinship recognition technology from research-to-reality. With clear problem statements, and an established measure for the practical significance, we compile a list of the SOTA scores and methods of the main tasks with emphasis on our large-scale FIW, which will be revisited in the next chapter with great detail.

In light of recent advances in our annual RFIW data challenge, along with the new task evaluations added as part of the most recent 2020 edition, we review the details for the different paradigms of kin-based problems in the visual domain as formulated for our large-scale, multi-task FIW database. We also look back at the existing datasets for visual kin-based problems motivated by different real-world scenarios (Table 4.1). In summary, our overarching goal is to pour the foundation for a deeper understanding of visual kinship recognition problem domain.

Even still, there are many unanswered questions that are high in potential for future efforts in machine vision research, like studies on familial and inheritance (*i.e.*, nature-based), and beyond. We take a glimpse at these promising next steps, while highlighting key challenges that we must overcome - both intrinsic to the image and inherent to the problem. We strongly urge there be attempts in research to establish cross-discipline studies— we are so ever ready to form such synergy.

The current chapter is of the form of a literature review in kinship recognition. There are a few related manuscripts that, together, span the material presented here.



**Figure 4.1: A decade of research in visual kinship recognition.** The timeline shows correlations between the data resources (*below timeline*) and citation metrics and events indicating the amount of research impact (*above timeline*). We built a pipeline to scrape the data needed for the plots above: (1) *Publish or Perish* [2] was installed on a Mac Book Pro to gather metadata for publications from various sources (*i.e.*, Google Scholar, Cross Ref, and Scopus) into a CSV file; (2) metadata in CSV was parsed into a BIB file using Python; (3) *Mendeley Reference Manager* was used to automatically detect duplicates while keeping as much information as possible by merging reference listings; (4) queried Google Scholar for all *Related Works* and *Cited By* using PyPi’s scholarly (<https://pypi.org/project/scholarly/>), which extended the paper-pile and increased the amount of metadata available from the richer metadata accessible using scholarly (*e.g.*, paper abstracts); (5) we clustered the documents by abstract via TF-IDF [3]. The clusters were high in recall, as true clusters were a majority of papers on kinship recognition in multimedia: this reduced the burden of manual inspection of hundreds of thousands to thousands. It is important to note that only citation metrics were considered, leaving out other factors of impact like the *number of times tweeted*, *Github stars*, and other indicators of impacting research.

**Table 4.1: Publicly available datasets for kinship recognition.** Each listed by the original name per reference. Kin-based image (or video) stats, which include the label types that support a specific evaluation metric and the respective SOTA score. URLs to the project page of each data resource are included. Abbreviations used for *Stats* are for the family count (**F**), face count (**f**), number of unique people (**P**), sample count (**S**), image count (**I**), video count (**V**), and multimedia (**MM**).

DB	Ref(s)	Stats	Label types	Metric, performance, SOTA	Web
CornellKin	[34]	• 150 <b>F</b> • 300 <b>S</b> • 300 <b>f</b>	parent-child	verification accuracy 94.4% [104]	<a href="http://chenlab.ece.cornell.edu/">chenlab.ece.cornell.edu</a>
UB Face	[105] [106]	• 200 <b>F</b> • 250 <b>P</b> • 600 <b>f</b> • 400 <b>S</b> • 1,736 (finger, 3D face, iris, DNA) <b>S</b> • 197 <b>I</b>	(young, old) parent-child	accuracy, 95.3% [104]	<a href="http://www.l.ece.neu.edu/yunfu/twins-day-dataset-2010-1015">www.l.ece.neu.edu/yunfu/twins-day-dataset-2010-1015</a>
Twins Day	[107]	• 184 <b>S</b> • 78 <b>P</b> • 78 <b>F</b> • 184 <b>f</b>	twin pairs	accuracy, 98.8% [107]	
SibFace	[108]	• 184 <b>S</b> • 78 <b>P</b> • 78 <b>F</b> • 184 <b>f</b>	siblings (brothers, sisters, mixed)	accuracy, 52.5% [109]	<a href="http://areeweb.polito.it/ricerca/">areeweb.polito.it/ricerca/</a>
UVa-NEMO Smile	[110]	• 162 <b>P</b> • 515 <b>V</b> • 512 <b>S</b>	7 relationships (core family)	accuracy, 88.16% [111]	<a href="https://www.uva-nemo.org/">https://www.uva-nemo.org/</a>
Family101	[9]	• 101 <b>F</b> • 607 <b>S</b>	family-tree structure	rank@10, 70.1% [29]	<a href="http://chenlab.ece.cornell.edu/">chenlab.ece.cornell.edu/</a>
KFW I + II	[112]	• 533 + 1,000 <b>P</b> • 1,066 + 2,900 <b>f</b>	parent-child; same + different photo	accuracy, 96.9% + 97.1% [104]	<a href="http://www.kinfacew.com/">http://www.kinfacew.com/</a>
TSKin	[113]	• 787 <b>F</b> • 2,589 <b>S</b>	(father & mother)-child	accuracy, 91.4% [114]	<a href="http://parnec.nuaa.edu.cn/TSKinFace">parnec.nuaa.edu.cn/TSKinFace</a>
FIW	[6, 7]	• 1,000 <b>F</b> • 33,000 <b>f</b> • 1-M <b>P</b> • 12,000 <b>S</b> • 13,000 <b>I</b>	large-scale; person-, family-, and image-level	accuracy, 78%; tri-subject accuracy, 79%; mAP 18% & rank@5 60% [7]	<a href="https://web.northeastern.edu/smilelab/fiw/">https://web.northeastern.edu/smilelab/fiw/</a>
KFWV	[115]	• 418 (video) <b>P</b> (100-500 <b>f</b> per video);	parent-child	accuracy, 61.8% [115]	<a href="https://www.kinfacew.com">https://www.kinfacew.com</a>
KIVI	[104]	• 503 (video) <b>S</b> • 503 <b>I</b>	7 relationships (core family)	accuracy, 83.2% [104]	<a href="http://fab-rubric.org">http://fab-rubric.org</a>
FIW-MM	[19]	FIW + 937 <b>MM</b>	FIW + multimedia for ≈ 200 <b>F</b> ( <i>i.e.</i> , video, audio, and contextual data)	EER, 89.8%; mAP, 0.24 [19]	<a href="https://web.northeastern.edu/smilelab/fiw/">https://web.northeastern.edu/smilelab/fiw/</a>

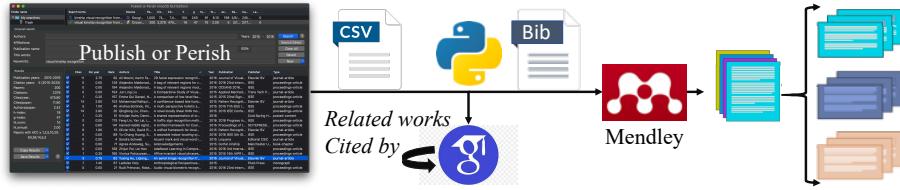


Figure 4.2: **Workflow to scrape publication metadata for Figure 4.1.** From *Publish or Perish* [2], we queried Scholar for *Related works* and *Cited by*, increasing the size of our list nearly 20-fold. Mendeley merged duplicates, while keeping as much information as possible. Applied NLP to cluster relevant documents.

Firstly, is the first survey on visual kinship recognition published in 2016, which was an extensive overview of the SOTA methods and data resources of the time [116]. The authors proposed future directions with great emphasis on the lack of labeled data both in sample counts and relationship label types. Hence, Wu *et al.* claimed that traditional metric-based solutions were inferior to deep learning models (*e.g.*, CNN), which has since shown to be true—evidence in the current SOTA, which we will introduce in the current chapter, with support and additional details provided in the chapters that proceed. Hence, the release of our large-scale FIW dataset to support modern-day, data-driven solutions was introduced at the end of that very same year [4].

In 2018m Georgopoulos *et al.* surveyed kinship and age in FR [117]. Although a comprehensive piece, kin-based problems ought to be surveyed independently. Nonetheless, prior knowledge of one could benefit the other; knowledge of various soft biometrics tends to complement and are beneficial (*e.g.*, gender and emotion). On the one hand, a survey on age or kinship should mention the other; however, the directed graphs and concepts of inheritance make kin-based studies worthy of surveying as an independent topic. Finally, looking at the problem from the view of understanding age, we make a similar claim—knowledge of kinship could most certainly help an age estimate. For instance, we have photos of a father at one or more known ages, while we are tasked to predict the age of the son. The knowledge available in the set of faces of known age is available as a prior and be modeled accordingly. Point is, we do not mean that the different attribute-based face understanding tasks ought to be treated as independent. Note, we do claim a study on the modeling and analysis of kin-based media should be surveyed alone. Still, as we cover later, age, gender, and variations of in other attributes do provide additional challenges in kin-based tasks.

More recently, Qin *et al.* surveyed kinship recognition methods as being founded on *a measure of kinship traits or statistical learning* [118]. Furthermore, the groups were characterized for

being *low* or *mid-level features*, *metric learning*, or *transfer learning*. The authors reported scores for several kin-based datasets. Additionally, *human* performance compared to machines was included. As part of their work, the authors proposed a standardized vision system based on four-steps to provide a generic, modular solution. To complement this, we define the problems consistently for the many kin-based tasks, and with details on the SOTA for each.

Most recently, we published an extensive survey on the topic of automatic kinship understanding [18]. In part, our survey was motivated by it being the 10 year anniversary since the first work in machine vision was proposed with several benchmarks and labeled data made available to the community for research purposes. On the other hand, our FIW dataset attracted lots of attention the past few years for solution, and we review the major-milestones in research of automatic visual kinship understanding over the first decade. The discussion is supported by a detailed illustration to assess the problem as it evolved over time, along with the public data supporting the progress, and with data statistics and web links of the source. Furthermore, we look at kinship recognition research that compares humans to machines, showing resemblance is detectable via the human eye (Section 4.2). Next, we introduce kin-based tasks by discussing the different problem statements (Section 4.3). Following this, we discuss experimental details for each of the tasks— summarize the protocols of the laboratory-based evaluations, including the data splits, metrics, and baseline results for each (Section 5.7). Then, we cover methodologies, both traditional and deep learning based for both discriminative and generative (Section 5.8).

Then, we discuss technical challenges preventing this technology from working reliably in real-world applications. Specifically, we cover the current limitations of SOTA— raise the discussion on a more broad perspective of the impact from kin-based technologies (*i.e.*, in our everyday lives). This is supported by a rigorous analysis on the edge cases and commonalities of falsely predicted samples. We highlight challenges posed by nature and the environment, and then shine a light on the inherent difficulties of obtaining sufficient data for kin-based problem (Section 7.2). This leads to the applications that line up with specific task-evaluations, both existing (*i.e.*, practically existing) and high in potential (*i.e.*, hypothetically possible). Emphasis is especially placed on the more robust models— typically, assuming we can improve the performance of the current SOTA (Section 7.3).

## 4.2 Background Information

The story of visual kinship recognition research can be told through the data. Therefore, we speak of the progress through the first decade from the perspective of the resources available

## CHAPTER 4. VISUAL KINSHIP RECOGNITION

(Figure 4.1), and it is shown that interest has been contingent on the amount and quality of labeled data. We end by discussing the data challenges, workshops, and tutorials used to motivate researchers.

### 4.2.1 The evolution of the problem

An increasing number of researchers have focused their attention on the problem of learning families in photos since the seminal paper was published in 2010 [34]. The research progress had the past decade coincided with the supporting labeled data released in part to it. Following Figure 4.1, we will next look back at the problem.

A trend observed in the progress in visual kinship recognition over the past decade is its correlation with the respective data resources released for public use. Critical points in the research stemmed from the respective problem statements supported by data labeled for the task. Hence, to review the problem statements and protocols as the problem evolved over time.

Fang *et al.* proposed training machinery to visually discriminate between *KIN* and *NON-KIN* using various facial cues [34]. Specifically, the authors demonstrated an ability to verify kinship given a face pair. To support this, they built and released the first facial image-based kinship database called Cornell Kin. Cornell Kin consisted of 150 face pairs of type *parent-child* from the web (*i.e.*, public figures, politicians, and other famous persons). Next came biometric data of twins collected at an annual event called *Twins’ Day* [107]. This effort yielded in a collection of 197 individuals of multiple modalities (*i.e.*, finger prints, 3D face scans, images of irises, and DNA samples) spanning multiple years (*i.e.*, from 2011 onward, each year new samples for subjects were added). Shao *et al.* then proposed UB Face made-up of 250 parent-child, each supported by three samples (*i.e.*, child, parent at a younger age, and parent at an older) [105]. The motivation for the pairs having a sample of each parent at a young and old age was directly spawned up from consideration for the difficulty imposed by large age gaps [119]. Soon thereafter came Family101 [9]—the first image collection with knowledge of family tree information, with 101 trees and multiple samples per subject. In 2014, Kin-wild I & II then provided a rich collection of 2,000 parent-child pairs [120]—will be discussed in the following section, along with Section 4.3.1, describing this database had significant impact for many expert researchers who proposed clever metric learning methods. Following this came the *Tri-Subject Kinship* (TSKIN) dataset, which structured the problem differently: given a parents-child pair (*i.e.*, both parents and a child), determine *KIN* or *NON-KIN*. Then, came FIW, which remains the largest kin-based image collection up to today. FIW is the main data used for experiments in this survey—more information provided in detail in the sections to come. Finally, and most recently,

## CHAPTER 4. VISUAL KINSHIP RECOGNITION

was the release of multimedia collections in support of kin-based tasks. First of these was released in 2017 - *Kin-Faces in the Wild* (KinFaceW) released video data for parent-child pairs, which then allowed for richer, dynamic models to be trained across video frames. Last year, in 2019, came the release of kin-based data that also leveraged audio media, *i.e.*, TALKing KINship (TALKIN). Lastly, FIW was extended with multimedia data added to over 200 of its 1,000 families [19]. Specs of the aforementioned data (*e.g.*, label types, SOTA, reference links) are in Table 4.1. Furthermore, the advancements in methodologies are later covered in detail.

To build Figure 4.1, *Publish and Perish* was used to acquire the paper-pile for the analysis (Fig 4.2). For this, a series of queries was executed, each using *visual kinship recognition* as the keywords: Google Scholar, limited to 1,000 search results per query, was run two times (*i.e.*, 2010-2020 and 2015-2020); Scopus was queried from 2010-2020, as only 48 items were found; CrossRef, with a limit of 200 items per search, was queried by year of publication (*i.e.*, 2010-11, 2011-12, ..., 2015-16, 2016, 2017, ..., 2020). Notice that the years were set such that fewest were expected the first year, more in the first half of the decade, and the most in the latter half. Many papers returned were not on automatic kinship recognition in visual media. However, using the TF-IDF representation, we were able to quickly filter out irrelevant papers by semi-supervised clustering (*i.e.*, side-information-based) cosine-similarity k-means [6] with labels assumed positive for the papers with keywords or titles that contain *visual kinship recognition*. A rise in the number of annual publications indicates an increase in interest of researchers; the impact on the research community, as a whole, nearly grows exponentially (*i.e.*, citation count). The incentive provided by data challenges, along with the increase in labeled data, influence the attention given to kin-based problems in multimedia (MM).

### 4.2.2 Humans recognizing kinship in photos

Several works in vision research evaluate the ability of humans to detect kinship given a face pair. In the seminar work, Fang *et al.* evaluated humans using a subset of their Cornell Kin dataset and found that, on average, humans were about 4% worse than the machinery: average human performance was 67.19%, with the top performance reaching 90% and the worst at just 50% [34]. The authors also found supporting results of a previous cognitive study on the perceivable similarity of offspring of different genders (*i.e.*, sons tend to be more recognizable as kin than daughters). Provided the time of this work- a time when it was still unclear whether *KIN* from *NON-KIN* signals are detectable via facial cues- this contribution was not only to compare humans and machines but

## CHAPTER 4. VISUAL KINSHIP RECOGNITION

was to get a sense if even possible for humans.

FIW was later used to evaluate humans at a greater scale, and with more diverse data [6]. Specifically, eleven relationship types, opposed to the four from Cornell Kin. Furthermore, pair counts were ten-times the prior. As described in the following subsection, the findings on FIW were that the machines outperform humans with the unconstrained data; additionally, and surprisingly, it was shown that it makes no difference, on average, if the human has prior knowledge of the relationship type in question.

More recently, a study focused on comparing human performance verifying kinship on a grander scale. Furthermore, the study presented extended coverage in the topic of human performance from the view of psychology. That is the work done by Lopez *et al.* [121]. Their evaluation displayed a face pair and prompted over 300 individuals connected through crowd sourcing services to answer: *Are these two people related (i.e., part of the same family)*, with possible responses set as Yes or No. Face pairs were made up of possible and negatives from both Kin-Wild [112] (*i.e.*, image set) and Uva Nemo Smile [110] (*i.e.*, video set) [122]. The machine again outperformed the human. Taking it a step further, Hettiachchi *et al.* analyzed the effects of gender and race, in both the human and data, finding that both genders are similar in ability to recognize kin, while both tend to perform best on same gender pairs (*e.g.*, brother-brother, mother-daughter, etc.) [123]. Furthermore, the authors validated previous findings in own-race bias for humans recognizing kinship.

### 4.2.2.1 Results

We assess both human evaluations via box plots (see Fig. 5.12).

In *Case 1*, the minimum scores across most categories are below random (*i.e.*, < %50). In response, we confirmed that no single person scored lowest in more than 1 of the 9 categories. Another observation is the distribution of averages, and its mean of 57.5%, had the smallest variance—no average below 50% or above 67.5%, which indicates that no single, or more than a few subjects, dominated the average scores for the better or the worse. Examining the pairs where errors were made, three conclusions can be made: (1) especially for relationship types spanning 1 or more generations (*i.e.*, parents and grandparents), the common pairs consistently marked incorrectly are cases where the face of the expected elder is at a younger age or the face of the descendant appears older (*e.g.*, grandfather in his thirties and grandson in his fifties); (2) different ethnic groups typically made common mistakes on face pairs of different backgrounds; (3) females often deviated from males on the mistakes made that are common and across different ethnic groups—varying females



Figure 4.3: **Samples used for human evaluation.** Each column displays pairs most commonly marked correctly and incorrectly, and in cases where the correct answers were true and false. Specifically, TP, FP, TN, and FN are displayed, respectively. Each of these pairs was properly classified by the fine-tuned CNN.

were always the top scorer, but never the same twice. Apparently, nature and nurture can play a role in humans' ability to do kinship verification as well. There are many interesting directions for future work (*e.g.*, even larger and more diverse subject pool, or samples with added semantics like full body views or entire photos with background context).

For *Case 2*, we evaluated humans' ability to recognize kinship in faces, but, this time, without specifying the relationship. From this, we were aiming to determine whether the relationship direction and face age impacted human responses. Overall, the mean values barely changed, however, the set of pairs commonly marked wrong did—relationship direction does seem to worsen human ability to recognize kinship when the direction of the relationship contradicts with the age appearance of face pairs; however, in cases without the age contradiction, knowledge of the relationship type helps humans to determine whether or not the face pairs are of that type (*i.e.*, even though the set of common pairs incorrectly classified changed, the overall mean did not, as the average fell between 57-58% in both cases). Figure 4.3 shows face pairs most commonly classified correctly or incorrectly considering both cases.

Quantitatively, human performers scored an average of 57.5%. This is comparable to hand-crafted features such as LBP and SIFT, but nearly 15% lower than our fine-tuned CNN (*i.e.*, the SphereFace CNN fine-tuned for this experiment scores 72.15%).

### 4.3 Visual Kinship Problems

As mentioned, there are several kin-based tasks, each defined by specific protocols to best help control the experiment while simulating the use-case. Before we introduce the experimental details in the proceeding section, let us first introduce the various views of kin-based problems by the

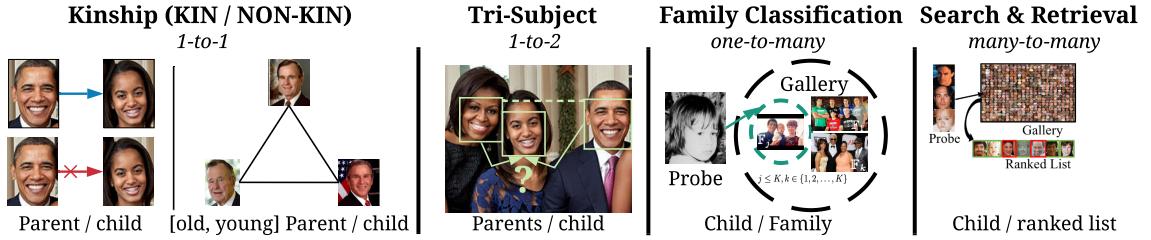


Figure 4.4: **Visual kin-based discriminate tasks for the FIW dataset.** Robinson *et al.* posed problems of verification (*i.e.*, one-to-one) [4] and family classification (*i.e.*, one-to-many) [5, 6], along with more recently supporting tri-subject verification (*i.e.*, one-to-two) and search & retrieval for “missing” children (*i.e.*, many-to-many) [7]: the FIW database supported the aforementioned tasks, while the annual data challenge RFIW was, and continues to be, motivated by promotional purposes. The most recent data challenge supported include three of the four shown here, as family classification was found to carry less potential for practical use-cases, while the others were done using three data splits disjoint in terms of family labels (Table 5.9, 5.11, 5.13). Protocols and benchmarks for each view are described in [7]. Best viewed electronically.

motivation and, thus, the problem statements for which they were found. Specifically, we review the discriminatory tasks (Figure 4.4), along with the generative.

### 4.3.1 Kinship verification

The goal of the most popular kin-based task is to determine whether a face pair are blood relatives (*i.e.*, *KIN* or *NON-KIN*). Scholarly findings in the fields of psychology and computer vision revealed that different types of kinship share different familial features. From this, the task has evolved into verification over a broad range of relationship types, like *parent-child* (*i.e.*, father-daughter (FD), father-son (FS), mother-daughter (MD), mother-son (MS)) or *siblings* (*i.e.*, BB, sister-sister (SS), brother-sister (SIBS)). Typically, we assume prior knowledge of the relationship type, both during training and testing. Hence, it is typical to train separate models or learn different metrics. Until the release of FIW [4, 5], small sample sets limited experiments. Thus, the larger data-pool of FIW resulted in larger-scale evaluations that better mimicked true distributions of diverse families globally. With it, also came additional relationship types that span multiple generations (*i.e.*, grandfather-granddaughter (GFGD), grandfather-grandson (GFGS), grandmother-granddaughter (GMGD), grandmother-grandson (GMGS)). FIW is made-up of 1,000 disjoint family trees of various

## CHAPTER 4. VISUAL KINSHIP RECOGNITION

structures (*i.e.*, the number of family members range from five to forty-four). Furthermore, subject nodes making up the trees typically contain multiple face samples— often samples that span over time, with face shots of most family members at different times in their lives. The families are split into five-folds with no overlap between folds. The trees are converted to pairs of various relationship types, with an average of five face samples per family member. The pairs present a variety of additional challenges, as, for instance, a GMGS pair may or may not be with faces of similar age. It could be an image of a younger grandmother, the GS as an adult, or maybe even as a GF himself.

Another flavor of kinship verification is best explained by the motivation behind UB Face: using knowledge of age as a prior and conditioned on whether or not *KIN* is true [105]. The idea was founded on analyzing the type of paired data frequently in the set of FP. Specifically, facial pairs of relatives separated by larger age gaps. Thus, based on perceived hard positives, the UB Face dataset provided a pair of images per parent— one at a younger and the other at an older age. In the end, Shao *et al.* supported their hypothesis experimentally by showing a pair of true *KIN* in *parent-child* relationships were closest when the parent was at a younger age. Then, Xia *et al.* used this to formally claim SOTA on UB Face by treating it as a transfer learning problem, with the target being that of the older parent and child, and the source being younger version of the respective parent and child [106]. Many have shown that paired data with greater age gaps are a challenge, and regardless of the level to which older children (*i.e.*, an elderly aged *child*) compares to older *parent*.

### 4.3.2 Family classification

Family classification, the problem where one family member is set aside, and all other members are used to model the classes (*i.e.*, family), is reviewed next. Hence, the task is to determine the family that the unknown subject belongs to, which is formulated as a closed-form, multi-class classification problem. This one-to-many problem is challenging, and only increases in difficulty with more families. The challenge stems from the large intra-class variations, which was revealed by a performance drop with an increasing number of families. Fang *et al.* [9] was first to demonstrate this on Family101. Specifically, the authors showed a drop in performance from ten-to-fifty families (increments of 10); opposite to this, the performance improved with one-to-four (increment of one) family member during training. Robinson *et al.* included 316 families originally [4], then 512 [6], and finally 564 [29]. After being supported as part of the RFIW annual data challenge the first three consecutive years (*i.e.*, 2017, 18, and 19), the overview of the latest RFIW mentioned the unrealistic setting of the problem, as families to employ on must be *a priori* knowledge (*i.e.*, unable

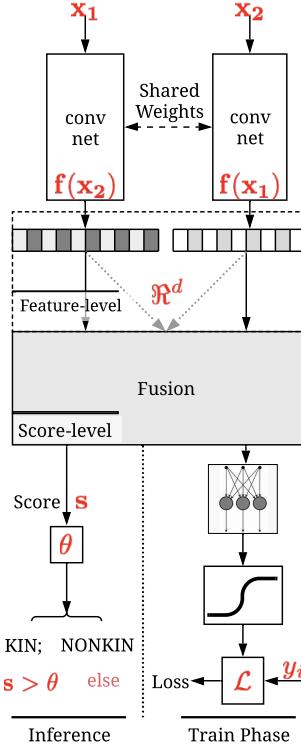


Figure 4.5: **Generic Siamese network.** Approaches tend to follow the Siamese model, differing in method of fusion. Specifically (from *top-to-bottom*), an image pair shot  $x_1$  and  $x_2$ .

to generalize well). Thus, Robinson *et al.* [7] omitted family classification from the latest challenge and substituted in two, more realistic views explained next. Nonetheless, as of the 2019 RFIW, 17.1% (accuracy) is SOTA [124].

### 4.3.3 Tri-subject verification

Tri-subject verification, first introduced in [113] (*i.e.*, TSKIN), focused on a slightly different view—predict whether a child is related to a pair of parents. In practice, this setting makes the most sense, as knowledge of one parent means the other can likely be easily inferred. Recently, as part of the RFIW data challenge, FIW was used to organize and benchmark tri-subject verification on scales much larger than ever before [7].

Table 4.2: **KinWild benchmarks.** Results for KinWild I and II.

	KinWild I					KinWild II				
	FD	FS	MD	MS	Avg.	FD	FS	MD	MS	Avg.
MNRML [125]	72.5	66.5	66.2	72.0	69.9	76.9	74.3	77.4	77.0	76.4
PDFL [126]	73.5	67.5	66.1	73.1	70.0	77.3	74.7	77.8	78.0	77.0
DMML [127]	74.5	69.5	69.5	75.5	72.3	78.5	76.5	78.5	79.5	78.3
multiviewSSL [128]	82.8	75.4	72.6	81.3	78.0	81.8	74.0	75.3	72.5	75.9
SSML [129]	81.7	75.3	71.4	77.9	79.6	82.4	78.6	79.8	77.9	79.7
SPML-P [130]	75.4	84.3	81.1	72.4	78.3	82.4	77.6	76.6	76.2	78.2
ELM [131]	70.0	64.2	73.0	77.2	71.1	78.6	73.6	81.0	79.6	78.2
KVRL-fcDBN [132]	<b>96.3</b>	<b>98.1</b>	<b>98.4</b>	<b>90.5</b>	<b>96.1</b>	<b>94.0</b>	<b>96.0</b>	<b>96.8</b>	<b>97.2</b>	<b>96.2</b>
MvGMML [133]	69.3	73.1	69.4	72.8	71.1	70.4	73.4	65.8	69.2	69.7
DDMML [134]	79.1	86.4	87.0	81.4	83.5	83.8	87.4	83.0	83.2	84.3
KML [135]	-	-	-	-	82.8	-	-	-	-	85.7
MSIDA+WCCN [136]	86.0	85.93	90.1	88.6	87.7	89.4	82.8	87.8	88.0	87.0

#### 4.3.4 Search and retrieval

This view, the most recent to be introduced [29], formulates the problem of missing (*i.e.*, unknown) children. A search *gallery* made up of all faces of FIW, but those of the single child held out as the *probes* for  $F$  families. Thus, the input is visual media of an individual, and the output is a ranked list of all subjects in the *gallery*. This *many-to-many* task is staged as a *closed set* problem. Thus, the number of TP varies for each subject, ranging anywhere from  $[1, k]$  relatives present in the *gallery*. In other words, there are always relatives present.

#### 4.3.5 Multi-modal data

Additional modalities (*e.g.*, video [104, 115], audio [137], multimedia [19]), although limited attempts and fairly new in literature, have proven quite effective. *KinFaceW Videos* (KFVW), spawned out of the same group as KinFaceW, meaning notable contributions by these authors at about the halfway point and towards the end of the decade (Figure 4.1). Wu *et al.* demonstrated that speech can be modeled to detect kinship [137]. In particularly, audio has shown promising,

but through minimal attempts. To better understand the patterns that allow for speech to work—whether that be jargon used, accents shared, or other acoustical features— we have seen that a kinship detection system can be improved with audio; however, an in-depth look at the model and the salient components of highly matched signals is subject to future work.

#### 4.3.6 Kin-based facial synthesis

Technology to post-process images (or even curate in real-time, *i.e.*, Snapchat filters) have grown popular in the modern-day main-stream. From this alone—kin-based face synthesis for entertainment and digital art is employable. As a concrete example, Snap Inc. introduced the ability to predict the offspring from a face pair of faces in their app mid-2019. Surely, a natural curiosity. Furthermore, studies support links between DNA and appearance [138], meaning it possible.

Another, nearly default use-case for synthesizing faces based on kinship is in law enforcement to predict the appearance of an unknown perpetrator provided knowledge of kin. Furthermore, missing family members (*e.g.*, a kidnapped child) with face images of years prior (*i.e.*, images of adolescence) could be used as prior knowledge, along with the appearances of family in their adulthood, to predict the face of that missing family member as an adult. Also, nature-based studies where latent variables control the appearance of an offspring in a manner that allows for the analyzer to explore. And, projecting further in time, presumably, is its place in genetics. If genetics allows for tweaking the fusion of male and female chromosomes to avoid face deformities of an offspring, the ability to visualize changes in appearance as a function of changes in latency, would likely be needed.

Some have attempted to predict the appearance of offspring in research (Section 5.8.3, while others seek a way to commercialize (Section 7.3): the former (*i.e.*, laboratory-style experimentation) and the latter (*i.e.*, applied in practical use-cases) are revisited in greater detail later.

# Chapter 5

## Families In the Wild (FIW)

### 5.1 Overview

Visual kinship recognition has an abundance of practical uses, such as issues of human trafficking and in missing children, problems from the modern-day refugee crises, and social media platforms. Use cases exist for academia as well, whether for machine vision (*e.g.*, reducing the search space in large-scale face retrieval) or a different field entirely (*e.g.*, historical & genealogical lineage studies). However, before the release of FIW in 2016 [4], no reliable system existed in practice. This is certainly not due to a lack of effort by researchers, as many works focused on kinship.

We identified two reasons that clearly slowed the rate at which visual kinship recognition technology evolved:

1. Data resources for visual kinship were too small to capture true data distributions.
2. Hidden factors in visual appearance shared by blood relatives are complex and less discriminant than in more conventional problems (*e.g.*, object classification or face identification).

A large image-set that properly represents families worldwide was needed. Firstly, the distribution one that also could meet the capacity of more complex, data-driven models (*i.e.*, deep learning). This was the primary motivation for us to build the first large-scale image database for kinship recognition, FIW. FIW is made-up of rich label information that captures the complex, hierarchical structures of 1,000 unique family trees. Families have an average of 13 family photos each (*i.e.*,  $>13,000$  family photos), and with family sizes that range from 3-38 members. Furthermore,

## CHAPTER 5. FAMILIES IN THE WILD (FIW)

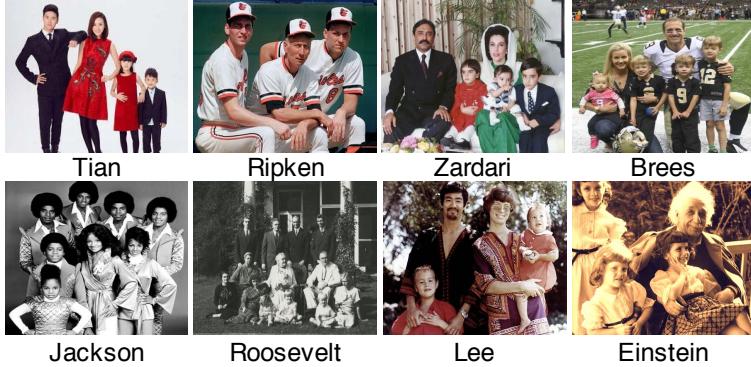


Figure 5.1: **Sample family photos from FIW.** Randomly picked 8 / 1,000 families.

most subjects have multiple samples across various ages (see Figure 5.1). FIW remains the **largest** and **most comprehensive** database of its kind.<sup>1</sup> Samples of a family from FIW is shown in Figure 5.9.

We proposed a multi-modal labeling scheme to minimize the amount of human input needed during the annotation process behind acquiring ground-truth for all the faces in FIW. As part of the process is a novel semi-supervised clustering method that proves to work effectively in practice (*i.e.*, generating label proposals for new data using existing labeled data as side information). Hence, after labeling about half of the data for half of the families via JAVA GUI, we could use these labels as the side information. Furthermore, we inferred face labels by comparing names in text metadata and the labeled faces and keeping the proposals of high confidence. A simple validation stage at the end takes minutes per a few set of families, opposed to hours. In the end, we show a significant reduction in manual labor and time spent on labeling new data.

Deep learning can now be applied to the problem, as we demonstrate on two benchmarks, kinship verification and family classification. We fine-tune deep models to improve all benchmarks, and provide details on the training procedure. We also measure human performance on verification and compare with benchmarks.

FIW was first introduced in [4], and later extended in [30], has constituted a number of contributions. These are listed as follows.

1. Added additional faces for verification and complete families for classification (Section 5.3).
2. Improved the labeling process with novel semi-supervised clustering method (Section 5.4).
3. Boosted baseline scores using up-to-date deep learning approaches (Section 5.5).

---

<sup>1</sup>Visit the FIW project page, <https://web.northeastern.edu/smilelab/fiw/>.

4. Obtained SOTA on smaller datasets via transferring CNN fine-tuned on FIW (Section 5.5.5).
5. Compared human performance with algorithms (Section 5.5.6).

## 5.2 Related Works

### 5.2.1 Related Databases

The story of visual kinship recognition begins in 2010, at which time the first kin-based image collection (*i.e.*, CornellKin) was made public [34]. CornellKin included 150 *parent-child* face pairs (*i.e.*, celebrities and their parents). Next, UB KinFace-I & II [119, 106, 139] were introduced to address a different view of kinship recognition— compare parent faces when young and old faces of parents were paired with a child, with a total of 600 face photos of 400 unique subjects (*i.e.*, celebrities and politicians). Then, KinWild I-II [125] was released and used in a 2015 FG Challenge [112], which too focused on *parent-child* pairs. Shortly thereafter, Family101 [9] was introduced as the first attempt of multi-class classification (*i.e.*, *one-to-many*) for kinship recognition. Thus, it is an organized set of structured families [9], including 206 sets of parents and their children (*i.e.*, *core families*) that make up 101 unique family trees. In 2015, TSKinFace [113] was built to support yet another view of kinship recognition, tri-subject verification, where both parents and a child were used— 513 Father/Mother-Daughter pairs and 502 Father/Mother-Son pairs (*i.e.*, *two-to-one* verification).

However, even after all these contributions, there existed no single resource that satisfied the concerns of insufficient data. A single resource with the features of previous works, but in a more complete and abundant manner, was the underlying vision for FIW. As shown in Tables 5.1 & 5.2, and discussed in later sections, FIW far exceeds others in terms of number of families, face pairs, and relationship types.

### 5.2.2 Automatic Kinship Recognition

As mentioned, Fang *et al.* first attempted kinship verification on *parent-child* face pairs [34]. They proposed selecting the 14 (of 44) most effective hand-crafted features. Following this, researchers recognized that a child’s face more closely resembles their parents at younger ages [119, 106, 139]. In response, they used transfer subspace learning methods that uses the younger faces of parents to help fill the appearance gap between their older faces and that of their children. To benchmark the KinWild dataset, Lu *et al.* proposed a metric learning method used in Euclidean space called Neighborhood Repulsed Metric Learning (NRML) and its multi-view counterpart (MNRML)

## CHAPTER 5. FAMILIES IN THE WILD (FIW)

that learns a common distance metric for multiple feature types [120]. Fang *et al.* focused on *one-to-many* (*i.e.*, family classification) by representing faces as a linear combination of sparse features (*i.e.*, feature selection via lasso) of 12 facial parts encoded via a learned dictionary [9].

Progress made in kinship recognition, along with release of varying task protocols, coincides with an increasing availability of structured and labeled data. Although there have been several significant contributions, none have overcome the challenges posed earlier.

### 5.2.3 Deep Kinship Recognition

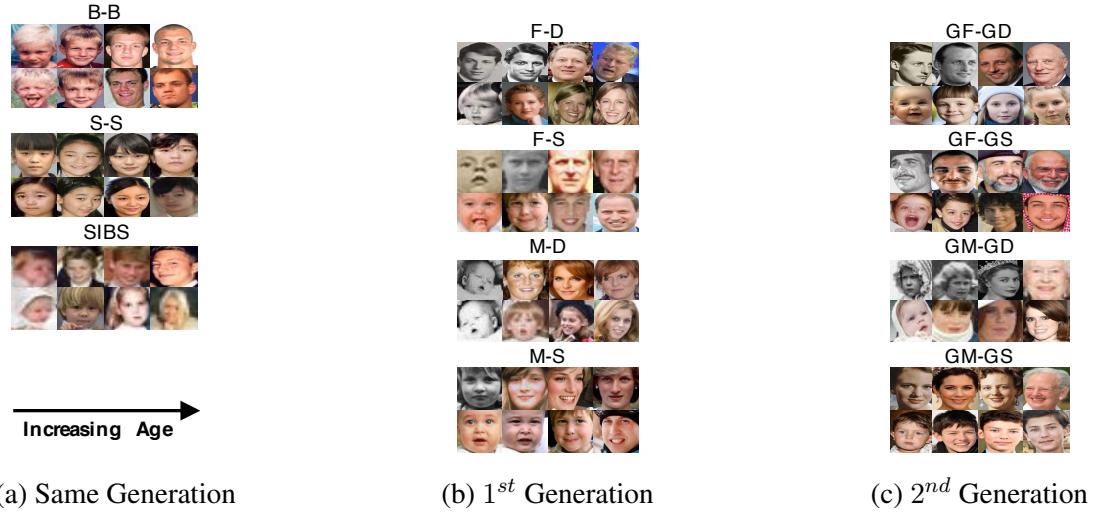
Since the AlexNet CNN [43] won the 2012 ImageNet Challenge [140], deep learning has achieved SOTA in a wide range of machine learning tasks. Central to this frenzy has been facial recognition [141, 46, 142]. In spite of this, there are only a few works that use deep learning for kinship recognition [143, 8, 144, 31].

Deep learning has yet to show an advantage for visual kinship recognition, with metric learning seeming more promising. As mentioned in a recent literature review [116], the reason for this stems from an insufficient amount of data. In this work, we include several benchmarks on FIW using deep learning, obtaining a clear advantage in both tasks.

### 5.2.4 Semi-Automatic Image Tagging & Data Exploration

Automatic image tagging was recently done by first labeling a small amount of the data, and then using it as side information to help guide the clustering process in a semi-supervised manner [145]. Following this, we take advantage of side information from labeled FIW.

Previous works used image captions, whether from Flickr or other sources of images tagged by users, to discover labels and annotate images in an automatic fashion [146]. Generally, methods mining text for image tags treat it as a problem of noisy labels [147, 148]. CASIA-WebFace [149], a large-scale dataset for facial recognition, successfully extended the scale of the renowned LFW [150]. By crawling the web, and leveraging knowledge from IMDB, multiple face samples for 10,000 unique subjects were collected. Although related in the sense of automatic labeling, these problems are very different from the one we present here. We aim to add more data to underrepresented families of the FIW database, and doing so by using the existing labels for each family as side information to guide our semi-supervised clustering method. We wish to maximize the number of labeled faces available to facilitate the clustering in order to generate label proposals. For this, we use the existing FIW labeled faces and the text metadata of the unlabeled data to automatically tag



**Figure 5.2: Samples of eleven pair types of FIW.** Each type is of a unique pair randomly selected from a set of diverse families, while four faces of each individual depict age variations.

faces using an iterative process governed by both visual and contextual evidence. As discussed in Section 5.5.4, our method consistently improves with increasing amounts of side information.

### 5.3 Families in the Wild (FIW) Database

We next cover the FIW database by first recalling the original FIW dataset and old labeling scheme [4]. Then, we describe the improved semi-automatic labeling process that enabled the collection to grow as large as it did. Finally, we compare the two.

### 5.3.1 FIW v0.1

Our goal for FIW was to collect about 10 family photos for 1,000 unique families and support with 2 types ground-truth labels, photo-level (*i.e.*, who and where in the image) and family-level (*i.e.*, all members and the relationships between them). Fig. 5.4 depicts the 2 label types. FIW was organized as follows: each family was assigned a unique ID (*i.e.*, FID), and pictures collected were also assigned a unique ID (*i.e.*, PID). Finally, members added were assigned their own unique ID (*i.e.*, MID). For instance,  $FID_1 \rightarrow MID_1$  in  $PID_1$  refers to the first member of the first family in the first photo collected. The order of IDs was arbitrary, as assignments were made in the order that the family, member, and photo were added. Before introducing the new and improved

## CHAPTER 5. FAMILIES IN THE WILD (FIW)

semi-automatic process, we briefly review the process used initially in [4], which involved 3 steps: (1) *Data Collection*, (2) *Data Labeling*, and (3) *Data Parsing*.

For *Data Collection*, a team of 8 students from different parts of the world, and with vast knowledge of famous persons, compiled a list of families with a primary focus on their place of origin (*i.e.*, an attempt to compile a diverse family list). Table 5.3 lists the ethnicity distributions of the 1,000 families. Note that this is not the exact distribution, as each family was counted once according to the *root* member for which the search was based (*i.e.*, not per member, but per family). For instance, for Spielberg’s family we consider just Stephen. Future work could entail adding more families from underrepresented ethnic groups, as the distribution still favors Caucasians.

For *Data Preparation*, we built a labeling tool to guide the process of generating the two label types. Annotators work through all family photos on a family-by-family basis, specifying who was in each photo by clicking member faces and choosing their names from a drop-down menu. Names, genders, and relationships for members were only entered on the first instance in an image—once added to the family then only must select their names upon clicking where in photo.

For *Data Parsing*, all family photos were detected using classic HOG features trained on top of a linear classifier using image pyramids and sliding windows via DLIB [151]. Faces were cropped and normalized as done in [64], and then resized to  $224 \times 224$ . Finally, the structure of the database was organized into a hierarchy of directories, FID→MID→Face-ID (*i.e.*, 1,000 folders, F0001-F1000, containing family labels and folders for MIDs with face samples of that member).

Even though it only took a small team to label 10,676 family photos and 1,000 families, the process relied heavily on human input. Plus, in the end, many families were not properly represented (*i.e.*, either too few members, face samples, or family photos). Thus, we aim to reduce the manual labor and overall time requirements to add additional data provided various amounts of labels existed for each (*i.e.*, 61 existing families and 4 replacement). We added replacement families (*i.e.*, newly added families) to make up for cases of overlapping families or an insufficient online presence when searching for photos (*i.e.*, unable to locate family photos for 2 of the under-represented families). Before we propose the semi-automatic labeling model, we first review the two benchmarks included in this work, along with the related statistics of each. We then present the new labeling process that enabled us to add additional data with far less manual labor and in just a fraction of the time.

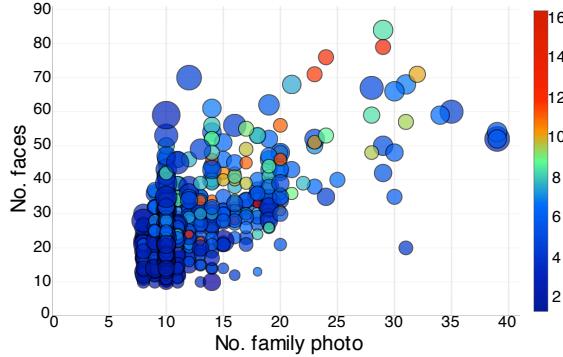


Figure 5.3: **Database statistics.** Horizontal and vertical axes represent counts for photos and faces per family, respectively. Bubble size and color represent counts for members and average faces per member, respectively.

### 5.3.2 Data Preparation

Due to the nature of the label structure, FIW can serve as a resource for various types of vision tasks. For this, we benchmark the two popular tracks, kinship verification and family classification. Next, we introduce both of these tasks and the means of preparing the data.

#### 5.3.2.1 Kinship verification

Kinship verification aims to determine whether two faces are blood relatives (*i.e.*, kin or non-kin). Prior research mainly focused on *parent-child* pairs (*i.e.*, father-daughter (F-D), father-son (F-S), mother-daughter (M-D), and mother-son (M-S)); some considered sibling pairs (*i.e.*, brother-brother (B-B), sister-sister (S-S), and brother-sister/mixed gender siblings (SIBS)). However, research in both psychology and computer vision revealed that different kin relations render different familial features, which motivated researchers to model different relationship types independently. With the existing image datasets used for kinship verification limited to, at most, 1,000 faces and typically only 4 relationship types, we believe such minimal data leads to overfitting and, hence, models that do not generalize well to unseen data captured *in the wild*. FIW currently supports 11 relationship types (see Figure 5.2), 4 being introduced to the research community for the first-time (*i.e.*, *grandparent-grandchild*) and, most importantly, each category contains many more pairs—418,000 face pairs in [4] has increased to 656,954 after extending FIW via the proposed semi-supervised approach.

The 11 relationship types provide a more accurate representation for real-world scenarios. As mentioned, FIW was structured such that the labels can be parsed for different types of tasks and

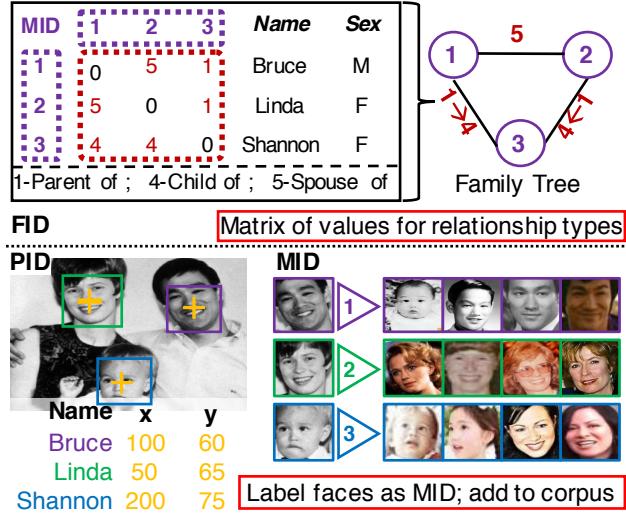


Figure 5.4: **Visual of the label types of FIW, Family-level (FID) and Photo-level (PID).** FID has individual family member (MID) and relationship information. PIDs contain information of MIDs + their locations in photos.

experiments, and additional kinship types can easily be inferred.

### 5.3.2.2 Family classification

Family classification aims to determine the family an unknown subject belongs to. Families are modeled using the faces of all but one family member, with the member left out used for testing. This *one-to-many* classification problem is a challenging problem that gets more challenging with more families. This is because families contain large intra-class variations that typically fool the feature extractors and classifiers, and each additional family further adds to the complexity of the problem. Additionally, and like conventional facial recognition, when the target is unconstrained faces *in the wild* [150] (*e.g.*, the variation in pose, illumination, expression, etc.), the problem continues to become more difficult. In [4], the experiment included only 316 families (*i.e.*, families with  $>5$  members). In this extended version, we now can include 524 families with the added data. We next present the process followed to extend FIW.

### 5.3.3 Extending FIW

The proposed semi-supervised model was used to generate label proposals, using existing and newly labeled data as side information for clustering— More side information consistently yields

## CHAPTER 5. FAMILIES IN THE WILD (FIW)

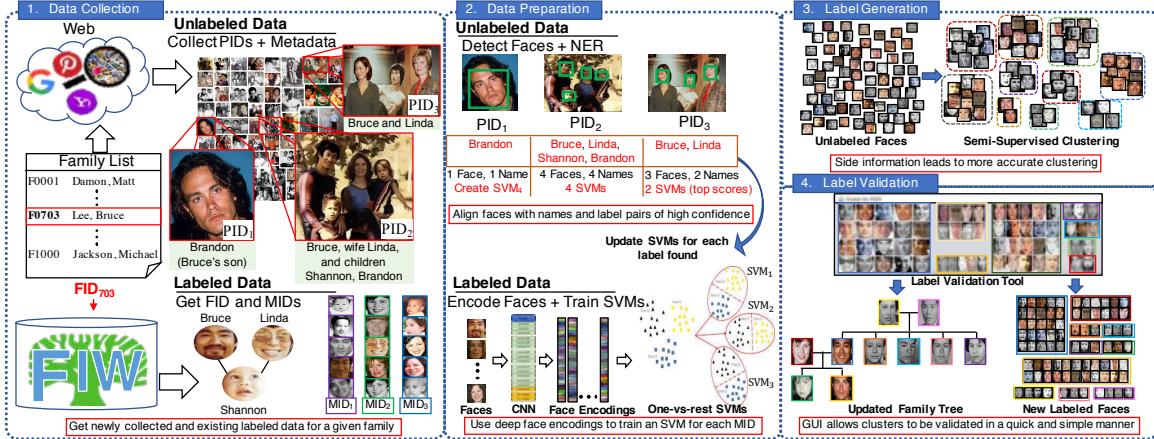


Figure 5.5: **Semi-automatic labeling pipeline.** *Data Collection.* Photos and text metadata were collected for underrepresented families in FIW and assigned unique IDs (*i.e.*, PIDs). Each new member requires at least 1 profile picture (*e.g.*, Brandon in *PID*<sub>1</sub>) to add to known labels. *Data Preparation.* With the existing FIW labels, we next aim to increase the amount, both in labeled faces and member labels, using multiple modalities—names in metadata and scores of SVMs were used to automatically label some unlabeled data—face-name pairs were assumed labeled for cases of high confidence. Starting from profile pictures (*i.e.*, 1 face, 1 name) and working towards less trivial scenarios (*e.g.*, 3 faces and 2 names, with 2 faces from 1 member at different ages, like in *PID*<sub>3</sub>). This step adds to the amount of side information used for clustering. *Label Generation.* Label proposals for remaining unlabeled faces were generated using the proposed semi-supervised clustering model that leverages labeled data as side information to better guide the process. *Label Validation.* A GUI designed to validate clusters and ensure clusters matched the proper labels.

better performance (see Section 5.5.4 & Fig. 5.11). Thus, we set out to maximize the amount of side information (*i.e.*, labeled faces) by inferring labels with high confidence by aligning faces and names from the unlabeled photos and metadata. Additionally, we modeled labeled data to discriminate between different family members in a photo. A single family was processed at a time to reduce both the search and label spaces. We aimed to discover labels using evidence from multiple modalities (*i.e.*, visual and contextual). This increased the amount of side information for clustering, and also the sample count modeled and used to label more faces. Resulting clusters were then set as ground truth upon human verification in *Step 4*.

We demonstrate the effectiveness of the new labeling scheme by comparing the number of user inputs (*i.e.*, mouse clicks and keystrokes) and overall time with the process followed in [4]. It

## CHAPTER 5. FAMILIES IN THE WILD (FIW)

Table 5.1: **Pairwise counts of FIW.** Notice FIW is first to provide Grandparent-Grandchild pairs. Table 5.2 further characterizes that data, and Figure 5.2 shows samples from it.

	siblings			parent-child				grandparent-grandchild				Total
	B-B	S-S	SIBS	F-D	F-S	M-D	M-S	GF-GD	GF-GS	GM-GD	GM-GS	
KinWild I [120]	0	0	0	134	156	127	116	0	0	0	0	533
KinWild II [120]	0	0	0	250	250	250	250	0	0	0	0	1,000
Sibling Face [108]	232	211	277	0	0	0	0	0	0	0	0	720
Group Face [109]	40	32	53	69	69	62	70	0	0	0	0	395
FIW(Ours) [4]	<b>103,724</b>	<b>39,978</b>	<b>73,506</b>	<b>92,088</b>	<b>129,846</b>	<b>82,160</b>	<b>112,618</b>	<b>7,078</b>	<b>4,830</b>	<b>6,512</b>	<b>4,614</b>	<b>656,954</b>

took just a few inputs and a few minutes on average per family, opposed to hundreds of inputs and several minutes to over an hour (see Table 5.2).

We next explain the improved multi-modal scheme made-up of 4 steps: (1) *Data Collection*, (2) *Data Preparation*, (3) *Label Generation*, and (4) *Label Validation*. The goal of (1) and (2) is to gather and increase the amount of side information available for (3), while (4) is to ensure correct labels for all new data. In other words, we set out to increase the labeled sample pool (*i.e.*, side information) by inferring labels for unlabeled faces, which adds to the set of training exemplars. The faces that were still unlabeled in (3) were clustered using all labels as side information. All newly added data is then verified by a human. The process is illustrated in Figure 5.5 and described next.

### 5.3.3.1 Step 1: Data Collection

The goal was to collect additional data for under-represented families of FIW (*i.e.*, families lacking in number of members, faces, and/or family photos). There were 65 families extended in total, with 1 family replaced due to a lack of available data and 3 merging together due to family overlap (*i.e.*, Catherine, Duchess of Cambridge, and her immediate family merged with the *Royal* family, as her and Prince William have 2 kids and, thus, bridging the two families). Several new labels and relationships resulted from this merge, with the *Royal* family growing from 29 to 38 members, which is now the largest tree of FIW. See Figure 5.3 and Table 5.3 for FIW statistics.

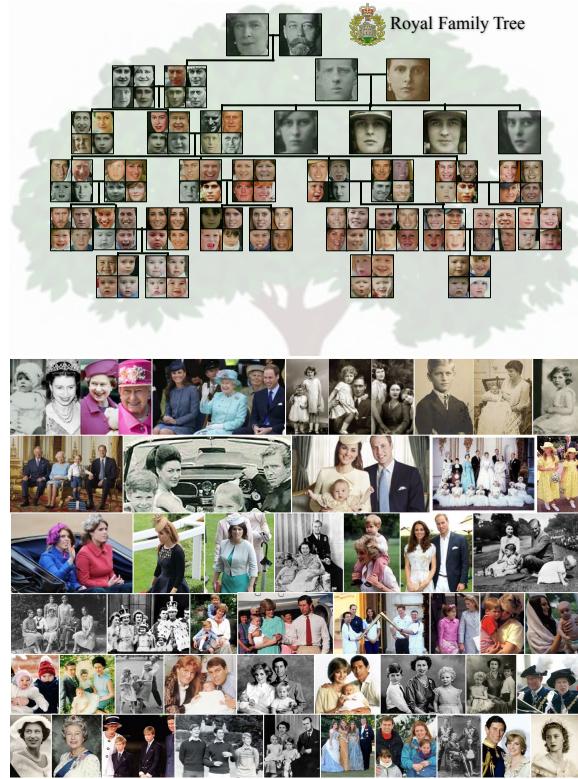


Figure 5.6: **Visualization depicting family structure and photos of the Royal Family.** There are several members in the tree (top) and many photos in total (bottom).

Preparing for *Step 2*, we set two requirements for the data: (1) rich text metadata describing subjects in each photo and (2) a profile photo per new member. Profile photos enabled label expansion for each new members (*i.e.*, single face and single name align with higher confidence).

### 5.3.3.2 Step 2: Data Preparation

The goal here was to maximize the amount of side information available for clustering in *Step 3*. Thus, we took advantage of both labeled (*i.e.*, faces & names) and unlabeled data (*i.e.*, detected faces & text metadata) to automatically infer labels for many unlabeled faces (see *Data Preparation* in Fig. 5.5). We next describe each component involved in this step.

**Text metadata** (*i.e.*, image captions) were collected for all photos in *Step 1*. With this, a list of names for each family was generated via a Name Entity Recognition (NER) classifier [152]. Then, a Look-Up-Table (LUT) of candidate names for each member was generated—*i.e.*, keys as member IDs (MIDs) and values as possible references to that member (*e.g.*, *Bruce* aka *Bruce Lee* aka *Brandon's*

## CHAPTER 5. FAMILIES IN THE WILD (FIW)



Figure 5.7: **Family photo montage.** Samples photos for 8 of 1,000 families in FIW.

CHAPTER 5. FAMILIES IN THE WILD (FIW)

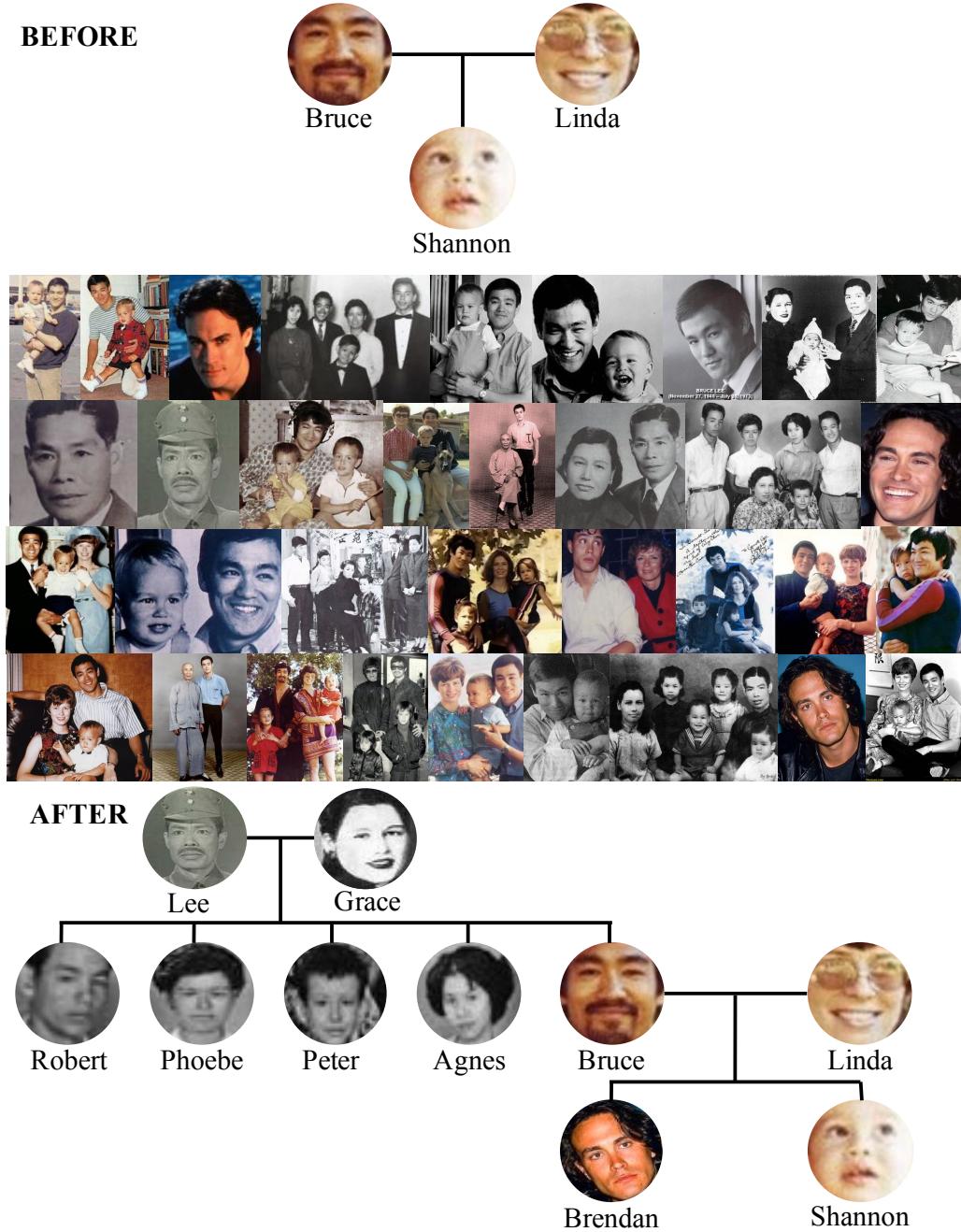


Figure 5.8: **Bruce Lee family tree before (top) and after (bottom) extension.** The photos in the middle were added to existing photos using our proposed semi-automatic labeling scheme. This increased both the samples per members and the total number (*i.e.*, from just 3 to 10 members).

## CHAPTER 5. FAMILIES IN THE WILD (FIW)

*father*). One challenge stemmed from name variations (*e.g.*, a person legally named *Joseph* might be called *Joe*); additionally, there were name titles (*e.g.*, *Queen Elizabeth II* might be called *Elizabeth II*, *Elizabeth*, or, in older photos, *Princess Elizabeth*). Additionally, nicknames posed additional challenges (*e.g.*, *Robert Gronkowski*, commonly referred to by his nickname *Gronk*). To address this, a LUT with detected references for each subject was first compiled, and then further augmented with additional tags (*e.g.*, adding titles and surnames). LUTs were later referenced to find evidence in the text metadata of a member’s presence in a photo.

**New MID**s found in profile photos (*e.g.*,  $PID_1$  in Figure 5.5)– when processing a family, each image that has a single face detected and just one name in its metadata was considered a profile photo. Profile photos were processed first. The name detected in the metadata was compared to all names for members stored in the LUTs. If there were no matches, the subject was then added as a new member in that family. A LUT of names was then generated for each new member, and the name of highest frequency (*i.e.*, number of detections in all metadata) recorded as the name corresponding to their assigned MID (*e.g.*,  $MID_6$  for the sixth member).

**Unlabeled and labeled faces** were encoded as 4,096D features from the  $fc_7$ -layer of the pre-trained VGG-Face CNN model [46]. *One-vs-rest* SVM models were trained for each member using labeled samples from all other members of that family as the negatives. Next, profile photos were processed (*i.e.*, 1 name and 1 face). Names that match an existing label were added to corresponding MID data pools, while mismatched names were added as a new MID with a LUT generated. This shows the benefit of including profile pictures for each new member, which makes it so all family members were known. It is important to note that SVMs were updated each time a new labeled face was added.

**Discovering labels** continues in a similar fashion, except now the SVMs play a more critical role. Now moving on to images with 2 faces and 2 names, the 2 SVMs of the respective members were used to classify the 2 faces. Provided high scores and no conflicts, labels were inferred. Cases with low confidence or conflicts were skipped, leaving those faces to be labeled via clustering. Next, photos with 3 faces and 3 names were processed, then 4 faces and 4 names, and so on and so fourth. After all one-to-one cases were processed, photos with a different number of names and faces were processed. For each photo, only SVMs that correspond to a LUT with matching names were used. Thus, justifying a requirement of *Step 1*– collect rich metadata in terms of specifying members present in photos.

Table 5.2: **Database counts and attributes.** Comparison of FIW with related datasets.

Dataset	No.	No.	No.	Age	Family
	Family	People	Faces	Varies	Trees
CornellKin[34]	150	300	300	✗	✗
UBKinFace[139, 106]	200	400	600	✓	✗
KFW-I[125]	✗	533	1,066	✗	✗
KFW-II[125]	✗	1,000	2,000	✗	✗
TSKinFace[113]	787	2,589	✗	✓	✓
Family101[9]	101	607	14,816	✓	✓
FIW [4]	<b>1,000</b>	<b>10,676</b>	<b>30,725</b>	✓	✓

Notice that some families benefited far more than others in this process. Nonetheless, roughly 25% of the 2,973 added faces were correctly labeled by this simple multi-modal process.

### 5.3.3.3 Step 3: Label Generation

Label proposals were generated for unlabeled faces using the proposed semi-supervised clustering method. To get the most out of our model we automatically labeled additional data in *Step 2*, while identifying all new members being added to each family. Hence, the number of members (*i.e.*,  $k$ ) was known for each family.

More details, including the objective function and solution, given in Section 5.4.

### 5.3.3.4 Step 4: Label Validation

Finally, clusters (*i.e.*, labels) were validated by a human. This was a three-part process: assign an MID to each cluster; validate each cluster, which was displayed in a grid of faces in the order of confidence score; specify gender and relationships of newly added members. As shown in Figure 5.5, a JAVA interface was designed to generate ground-truth for new data with just a few clicks of the mouse and minimal time per family. The inputs were cluster assignments for a family,

 Table 5.3: **Ethnicity distribution for the 1,000 of FIW.** *Mix* families contain  $>2$  ethnicity (*e.g.*, Bruce (*Asian*) and Linda (*Caucasian*) Lee with 2 children).

Caucasian	Spanish/Latino	Asian	African/AA	Arabic	Mix
64%	10.7%	9.1%	8.2%	2.0%	6.0%

CHAPTER 5. FAMILIES IN THE WILD (FIW)



Figure 5.9: **Sample family of FIW [6].** Faces and relationships of the American Football family, the Gronkowski's (*Top*). The montage shows less than half of all photos for respective family. Photo types are various, spanning profile faces (*top*) to images of different subgroups of family members. Furthermore, samples capture different times of life. Note, crops were made to fit montage (*Bottom*).

## CHAPTER 5. FAMILIES IN THE WILD (FIW)

Table 5.4: **Speedup analysis.** Previous (white) versus new (shaded) labeling processes compared in terms of inputs (keyboard and mouse clicks) and time (hours:minutes:seconds).

	Bruce Lee	Michael Jordan	John Malone	Craig Mccaw	Marco Reus	British Royal	Michael Jackson	Total
Inputs (count)	551	97	153	178	35	1,838	1,272	4,124
<b>Inputs (count)</b>	<b>12</b>	<b>6</b>	<b>10</b>	<b>15</b>	<b>7</b>	<b>21</b>	<b>24</b>	<b>95</b>
Time (h:m:s)	0:15:08	0:5:31	0:5:18	0:6:16	0:4:24	1:25:23	0:44:52	2:46:52
<b>Time (h:m:s)</b>	<b>0:1:11</b>	<b>0:0:31</b>	<b>0:1:05</b>	<b>0:0:56</b>	<b>0:0:31</b>	<b>0:6:44</b>	<b>0:7:13</b>	<b>0:18:11</b>

with faces listed in order of confidence (*i.e.*, cosine distance from centroid). MIDs were assigned in *Step 2* (*i.e.*, inferred from text, SVM scores, or both), which must also be validated. The outputs were labels for each PID and an updated relationship matrix (Figure 5.4).

### 5.3.3.5 Discussion

Seven families of various sizes were used to compare the old [4] and proposed labeling schemes—old scheme took 4,124 inputs in about 2.75 hours, and just 95 inputs in about 18.1 minutes via the new (see Table 5.4). Collecting and labeling the data for the extended FIW was done by a single person in days; it initially took a small team several months with the old scheme. Thus, demonstrating a significant savings in manual labor and time (note that greater amounts of data was originally collected, however, relative savings in time and manual labor clearly yields from process used in this extended version). A possible future direction is to use this scheme to extend families of FIW with video data. Another possibility is to use this method to extend the number of families, which, if on the order of thousands or more, then automating *Step 1* could further reduce savings (*i.e.*, web scrape for family information (*e.g.*, *Wiki*) and photos (*e.g.*, *Google*, *Bing*, etc.)).

## 5.4 Semi-Supervised Face Clustering

Labeling is a human-necessary and expensive task to benchmark data sets. Here we aim to accelerate the process by using some labeled data in advance. In this part, we demonstrate a novel semi-supervised clustering for labeling. Let  $X = \{x_i\}$  be the data matrix with  $n$  instances and  $m$  features and  $S$  be a  $n' \times K'$  side information matrix, which denotes  $n'$  labeled data instances into  $K'$  classes. Our goal was to make use of  $S$  to guide the remaining data into  $K$  classes, with  $K' \leq K$ .

### 5.4.1 Objective Function

Inspired by our previous work [145, 153], a partition level constraint is used to make the learnt partition agree with partial human labels as much as possible. To demonstrate the effectiveness of our labeling mode, K-means with cosine similarity is employed as the core clustering method to handle high-dimensional data due to its high efficiency and robustness. Our objective function is

$$\min \sum_{k=1}^K \sum_{x_i \in \mathcal{C}_k} f_{cos}(x_i, m_k) + \lambda U_c(S, H \otimes S), \quad (5.1)$$

where  $f_{cos}$  is the cosine similarity,  $H$  is the final partition,  $H_S = H \otimes S$  is part of  $H$  which the instances are also in the side information  $S$ ,  $m_k$  is the centroid of  $\mathcal{C}_k$ ,  $U_c$  is the well-known Categorical Utility Function [154] and  $\lambda$  is the trade-off parameter.

To better understand the last term in Eq. 5.1, we give the detailed calculation of  $U_c$ . Given two partitions  $S$  and  $H_S$  containing  $K'$  and  $K$  clusters, respectively. Let  $n_{kj}^{(S)}$  denote the number of data objects belonging to both cluster  $C_j^{(S)}$  in  $S$  and cluster  $C_k$  in  $H_S$ ,  $n_{k+} = \sum_{j=1}^{K'} n_{kj}^{(S)}$ , and  $n_{+j}^{(S)} = \sum_{k=1}^K n_{kj}^{(S)}$ ,  $1 \leq j \leq K'$ ,  $1 \leq k \leq K$ . Let  $p_{kj}^{(S)} = n_{kj}^{(S)} / n'$ ,  $p_{k+} = n_{k+} / n'$ , and  $p_{+j}^{(S)} = n_{+j}^{(S)} / n'$ . We then have a normalized contingency matrix (NCM), based on which a wide range of utility functions can be accordingly defined. For instance, the widely used category utility function can be computed as follows:

$$U_c(H_S, S) = \sum_{k=1}^K p_{k+} \sum_{j=1}^{K'} (\frac{p_{kj}^{(S)}}{p_{k+}})^2 - \sum_{j=1}^{K'} (p_{+j}^{(S)})^2. \quad (5.2)$$

It is worthy to note that  $U_c$  measures the similarity of two partitions, rather than two instances. The larger value of  $U_c$  indicates the higher similarity.

### 5.4.2 Solution

We notice that the first term in Eq. 5.1 is the standard K-means with cosine similarity. Could we still apply K-means optimization to solve the problem in Eq. 5.1? The answer is yes! Due to our previous work [155], we provide a new insight of  $U_c$  by the following lemma.

**Lemma 1.** *Given a fixed partition  $S$ , we have*

$$U_c(H_S, S) = -\|S - H_S G\|_F^2 + \text{constant}, \quad (5.3)$$

where  $G$  is the centroid matrix of  $S$  according to  $H_S$ .

By the above lemma, the second term in Eq. 5.1 can also be transformed into a K-means problem with squared Euclidean distance. Then a K-means-like algorithm can be used on the augmented matrix with modified distance function and centroid update rule for the final partition.

First an augmented matrix  $D$  is introduced as follows.

$$D = \begin{bmatrix} X_S & S \\ X_T & 0 \end{bmatrix} \text{ with } X = \begin{bmatrix} X_S \\ X_T \end{bmatrix}, \quad (5.4)$$

where  $d_i$  is the  $i^{th}$  row of  $D$ , which has of two parts,  $d_i^{(1)}$  and  $d_i^{(2)}$  (*i.e.*,  $d_i^{(1)} = (d_{i,1}, \dots, d_{i,d_m})$  presents the feature space and  $d_i^{(2)} = (d_{i,d_m+1}, \dots, d_{i,d_m+K'})$  denotes the label space). Zeros in  $D$  are the artificial elements, rather than the true values so that all zeros contribute to the computation of the distance and centroids, which inevitably interfere with the cluster structure. To make the zeros in  $D$  not involved in the calculation, we give the new update rule for the centroids of  $D$ . Let  $m_k = (m_k^{(1)}, m_k^{(2)})$  be the  $k^{th}$  centroid  $\mathcal{C}_k$  of  $D$ , we modify the computation of centroids as follows.

$$m_k^{(1)} = \frac{\sum_{d_i \in \mathcal{C}_k} d_i^{(1)}}{|\mathcal{C}_k|}, \quad m_k^{(2)} = \frac{\sum_{d_i \in \mathcal{C}_k} d_i^{(2)}}{|\mathcal{C}_k \cap X_S|}. \quad (5.5)$$

and the distance function is also adjusted as

$$f(d_i, m_k) = f_{cos}(d_i^{(1)}, m_k^{(1)}) + \mathbf{1}(d_i \in S) f_{sqE}(d_i^{(2)}, m_k^{(2)}), \quad (5.6)$$

where  $\mathbf{1}$  returns 1 if the condition is satisfied, otherwise 0.

The correctness and convergence of the modified K-means is similar to one in [145].

## 5.5 Experiments

We conduct the following experiments: benchmark kinship verification and family classification; evaluate the proposed semi-supervised clustering method at the core of the new labeling scheme; fine-tune CNNs using FIW and evaluate on KinWild I & II (*i.e.*, transfer-learning); measure human performance on kinship verification and compare to top scoring algorithms.

The subsequent subsections are organized as follows. First, we review the visual features, metric learning methods, and deep learning common in all experiments. Then, we dive into the experiments mentioned above. We introduce each independently, but with the same structure: experimental settings, experiment-specific training philosophy, and then the results and analysis.

Table 5.5: **Averaged verification accuracy scores (%) for 5-fold experiment on FIW.** Note that there was no family overlap between folds.

Method	siblings			parent-child				grandparent-grandchild				Acc. ± Std.
	B-B	S-S	SIBS	F-D	F-S	M-D	M-S	GF-GD	GF-GS	GM-GD	GM-GS	
LBP [156]	55.52	57.49	55.39	55.05	53.77	55.69	54.65	55.79	55.92	54.00	55.36	55.33 ± 1.01
SIFT [157]	57.86	59.34	56.91	56.37	56.24	55.05	56.45	57.25	55.35	57.29	56.74	56.80 ± 1.17
ResNet-22 [158]	65.57	69.65	60.12	59.45	60.27	61.45	59.37	55.37	58.15	59.74	59.70	61.34 ± 3.81
VGG-Face [46]	69.67	75.35	66.52	64.25	63.85	66.43	62.80	62.06	63.79	57.40	61.64	64.89 ± 4.68
+ITML [159]	57.15	61.61	56.98	58.07	54.73	57.26	59.09	62.52	59.60	62.08	59.92	59.00 ± 2.44
+LPP [160]	67.61	66.22	71.01	62.54	61.39	65.04	63.54	63.50	59.96	60.00	63.53	64.03 ± 3.32
+LMNN [161]	67.11	68.33	66.88	65.66	67.08	68.07	66.16	61.90	60.44	63.68	60.15	65.04 ± 3.00
+GmDAE [162]	68.05	68.55	67.33	66.53	68.30	68.15	66.71	62.10	63.93	63.84	63.10	66.05 ± 2.36
+DLML [163]	68.03	68.87	67.97	65.96	68.00	68.51	67.21	62.90	63.96	63.11	63.55	66.19 ± 2.36
+mDML [31]	69.10	70.15	68.11	67.90	66.24	70.39	67.40	65.20	<b>66.78</b>	63.11	63.45	67.07 ± 2.44
ResNet+CF [5]	69.88	69.54	69.54	68.15	67.73	71.09	68.63	<b>66.37</b>	66.45	<b>64.81</b>	64.39	67.87 ± 2.15
SphereFace[48]	<b>71.94</b>	<b>77.30</b>	<b>70.23</b>	<b>69.25</b>	<b>68.50</b>	<b>71.81</b>	<b>69.49</b>	66.07	66.36	64.58	<b>65.40</b>	<b>69.18</b> ± 3.68

## 5.5.1 Experimental Setting

For the sake of organization, all low-level features and metric learning approaches used throughout are listed and described in this section. Most are in two or more experiments, however, even those used for verification, for example, are still treated as common information, and thus is described alongside other items of preliminary information. Following traditional “shallow” methods, we review specifications of the pre-trained CNNs used as off-the-shelf face encoder.

### 5.5.1.1 Feature Representations

Detected and aligned faces were normalized and encoded using low-level and CNN-based features. We next describe the descriptors used in this work— SIFT, LBP, pre-trained VGG-Face and ResNet CNNs— each having been widely used in visual kinship and facial recognition problems.

**SIFT** [164] is amongst the most widely used feature type in object and face recognition. Here we follow the settings of [125]: resize images to  $64 \times 64$ , then extract features from  $16 \times 16$  blocks with a

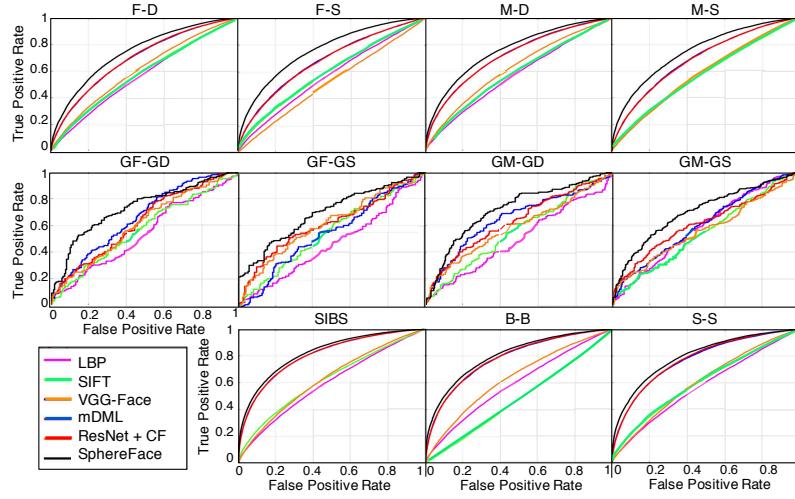


Figure 5.10: **Relationship type specific ROC curves.** Notice the fine-tuned Sphereface dominates, while the sample counts for the *grandparent-grandchild* were less as indicative of the jagged curves.

stride of 8 (*i.e.*, 49 blocks that yields  $128 \times 49 = 6,272$ D face feature).

**LBP** [156] is renown for its effectiveness in tasks such as texture analysis and face recognition. We again follow the settings of [125]: resize images to  $64 \times 64$ , divide into  $16 \times 16$  non-overlapping blocks, and use a radius of 2 and sampling number of 8. Each block was represented as a 256D histograms (*i.e.*,  $256 \times 16 = 4,096$ D face encoding).

**VGG-Face** [46], a pre-trained CNN with the topology of VGG-16: made-up of small convolutional kernels (*i.e.*,  $3 \times 3$ ) with a convolutional stride of 1 pixel. VGG-Face was trained on  $\approx 2.6M$  face images of 2,622 different celebrities. VGG has worked well on various face databases— 97.3% in accuracy on *YouTube Faces* [165]; 98.95% accuracy on *Labeled Faces in the Wild* [166]. By removing the top two layers— softmax and last fully-connected layer (aka  $fc_8$ -layer or  $fc_8$ )— the CNN can be used as an *off-the-shelf* face encoder [32]. Thus, models get trained on an auxiliary resource and employed on target data. Here, we fed faces through to the  $fc_7$ -layer (aka  $fc_7$ ), yielding a 4,096D face encoding.

**ResNet-22** [158] is a 22-layer residual CNN trained on CASIA-Webface [149]. ResNet-22 has a different network topology than VGG (*i.e.*, more layers made possible via skipping connections in residual blocks to ensure that the signal stays intact by superimposing an identity tensor). Faces were fed through to layer  $fc_5$  (512D encoding).

### 5.5.1.2 Metric Learning

Metric learning has been commonly used and, frequently, designed for kinship problems. Four metric learning and graph embedding methods used previously for face-based problems include: Information theoretic metric learning (ITML) [159], Discriminative Low-rank Metric Learning (DLML) [163], Locality Preserving Projections (LPP) [160], and Large Margin Nearest Neighbor (LMNN) [161].

### 5.5.1.3 Deep Learning

**Fine-Tuned CNNs.** Centerface (CF) [158] loss enhances the discriminative power of deeply learned features by adding a supervision signal to reduce the intra-class variations. SphereFace uses an angular softmax loss, and has most recently claimed state-of-the-art in facial recognition [48]. We fine-tune both these CNNs on FIW.

Additionally, we include two state-of-the-art methods based on AE, graph regularized marginalized Stacked AE (GmDAE) [162], and marginalized denoising AE (DAE) based metric learning (mDML) [31].

## 5.5.2 Kinship Verification

Kinship verification is a binary classification problem (*i.e.*, *true* or *false*, aka *kin* or *non-kin*, respectfully). It is the *one-to-one* view of kinship recognition, which we explain next.

### 5.5.2.1 Experimental Setting

The protocol followed is conventional in face-based tasks: 5-fold cross validation with no family-overlap between folds. There are 11 relationship types evaluated (summarized in Table 5.1).

For each pair type, we added negative (*i.e.*, *non-kin*) pairs to the 5-folds— we randomly mismatched pairs in each fold until the number of negative and positive pairs are the same in each fold (*i.e.*, negative pairs are added at random until it makes up 50% of the respective fold). Thus, the total number of positive and negative labels are equivalent.

For this task we included each feature, metric learning approach, and deep learning model listed above. We then fine-tuned the pre-trained CNN models on the FIW dataset, which is described in detail in the next subsection. To compare features, we computed cosine similarity between each pair, which was then compared to a threshold to classify each pair as either *kin* or *non-kin*.

## CHAPTER 5. FAMILIES IN THE WILD (FIW)

Verification accuracy (*i.e.*, average of 5-folds) and receiver operating characteristic (ROC) curves were used to evaluate.

### 5.5.2.2 Training Philosophy

For ResNet-22 + CF, we fine-tuned the Centerface model on our FIW data. Training was done using four Titan X GPUs with a batch size of 256. The learning rate was initially set to 0.01, then drops to 0.001 and 0.0001 at the 800 and 1,200 iterations, respectively. Training was complete after 1,600 iterations. The weight decay was set to 0.0005. For SphereFace [48], the settings are similar to ResNet-22+CF (*i.e.*, same batch size, learning rate, weight decay, and number of iterations), and with the angular margin set to 4.

### 5.5.2.3 Results

As listed in Table 5.5, *sibling* pair types tended to score the highest, followed by *parent-child* types, and then *grandparent-grandchild*. Thus, the wider the generational gap, the wider between appearances of faces.

SphereFace, which was fine-tuned on FIW, outperformed other benchmarks with an average accuracy of 69.18%, which is 1.31% and 2.11% better than ResNet-22+CF and mDML, respectively, which were the top scoring methods prior to the recent release of SphereFace. Also, out of the pre-trained CNNs, VGG-Face scored 3.55% higher than ResNet-22, and both outperformed the low-level features (*i.e.*, LBP & SIFT). From such, encodings from VGG-Face were used as features for the metric learning and AE methods. Besides LMNN and DLML, which improved score by 0.15% and 1.30%, the other metric learning methods actually worsened the performance of the descriptors extracted from the pre-trained VGG-Face CNN. This infers that faces encoded via VGG-Face are more discriminative when used *off-the-shelf* than when metrics are learned on top.

We show a significant boost in performance when fine-tuning CNNs on FIW data— all features from CNNs outperform the conventional shallow methods. The results show that the deep learning models better encode the complex representation needed to discriminate between *kin/non-kin* (see Figure 5.10). An improvement to these benchmarks, perhaps via a deep network designed specifically for this task, is certainly a direction for future work.

Table 5.6: **Family classification results.** Accuracy scores (%) using 564 families.

Run ID	Network(s)	Acc.
Run-1	VGG-Face, $fc_7$ (4,096D)+one-vs-rest SVMs	3.04
Run-2	VGG-Face, replaced softmax (564D)+fine-tuned	10.42
Run-3	ResNet-22 + softmax (564D)	14.17
Run-4	SphereFace (564D)	13.86
Run-5	ResNet-22 + CF (512D) + softmax (564D)	<b>16.18</b>

### 5.5.3 Family Classification

Family classification is a *one-to-many* problem. The goal is to determine which family an unseen subject came from. In other words, a set of families with a missing member to the model is provided. Then, the missing (*i.e.*, unseen) members get classified as being from one of the families (*i.e.*, closed form, as we currently assume that all members at test time belong to one of the families modeled during training). We next review some details for this task.

#### 5.5.3.1 Experimental Setting

Data from 564 families leaving a different single member out in each fold for testing, while data from all the other members was used for training (*i.e.*, leave-one-out w.r.t. family members). Families with at least 5 members were used. Thus, the data was split into 5-folds with no family overlap between folds (*i.e.*, a minimum of 4 family members for training and 1 for testing). Each fold contained roughly 2,700 images—about that many faces used to test each split, while the rest, about 12,800 faces, were used for training (*i.e.*, a total of 13,420 images).

#### 5.5.3.2 Training Philosophy

VGG-Face and ResNet-22 CNNs were fine-tuned on FIW by replacing the loss layers of the pre-trained CNNs with a softmax loss to predict the 564 family classes. There were a few differences: VGG-16 used a fixed learning rate of 0.0001, a batch size of 128, and trained for 800 iterations on one Titan X GPU; ResNet-22 used the same batch size and number of iterations, but with a larger learning rate 0.001, which was fixed too. For ResNet-22 + CF and SphereFace, we followed the same training process used for verification.

### 5.5.3.3 Results

We report the accuracy scores for five runs (see Table 5.6). As shown, the top-1 accuracy for modeling *one-vs-rest* linear SVMs on top of deep VGG-Face features was just 3.04%. Then, by replacing the softmax layer to target the number of families (*i.e.*, 564), and fine-tuning on FIW, the top-1 accuracy was improved (*i.e.*, +7.38% to 10.42%). ResNet-22, also fine-tuned by replacing softmax layer, showed the second to highest accuracy with 14.17%, which outscored the top performing CNN on verification (*i.e.*, SphereFace). The top performance was obtained with the fine-tuned ResNet-22 using Centerface (CF) loss with 16.18%.

### 5.5.4 Proposed Semi-Supervised Clustering

To demonstrate the effectiveness of our semi-supervised model, we cluster FIW data using various amounts of *family-level* labels as side information. We simulate two settings for evaluation—all data and just unlabeled data—shown as bold and dotted lines, respectively (see Fig. 5.11).

#### 5.5.4.1 Experimental Setting

We used 23,979 faces from 996 family classes. Faces were encoded using a pre-trained VGG-Face (*i.e.*,  $fc_7$ ). We varied the ratio of unlabeled data to side information across the horizontal axis up to 50% percent of labeled clusters, while the y-axis denotes the clustering performance on the rest of the unlabeled data by NMI. We compared to a pair-wise constrained clustering method, LCVQE [167], which is also a K-means-based constrained clustering method and transforms the partition level side information into ‘must-link’ and ‘cannot-link’ constraints. We used K-means as a baseline (*i.e.*, no side-information).

#### 5.5.4.2 Results

Fig. 5.11 shows a clear boost in performance for our method with more side information. Even on the unlabeled data, our method outperforms K-means, further validating the effectiveness of our method for semi-automatic labeling tasks. For LCVQE, the pair-wise constraints make the cluster structure unpredictable, vulnerable to deviate from the true one, and, thus, perform worse than the baseline. This shows that imposing hard constraints on side information, like ‘must-link’ and ‘cannot-link’, may even damper results. On the contrary, our model leverages the side information to only improve when more is added.

Table 5.7: **Transfer learning experiment.** Accuracy (%) for KinWild I & II. CNN fine-tuned on FIW top scorer. Note that these results were up-to-date when journal (*i.e.*, [6]) was released, but is no longer. See Table 4.2 for most up-to-date scores.

Method	KinWild-I					KinWild-II				
	FD	FS	MD	MS	Avg.	FD	FS	MD	MS	Avg.
LBP [156]	72.8	79.5	71.7	68.1	73.0	70.8	78.4	69.0	73.2	72.9
SIFT [157]	73.9	81.4	76.4	71.1	75.7	72.2	78.8	82.2	79.6	78.2
NRML (LBP) [125]	81.4	69.8	67.2	72.9	72.8	79.2	71.6	72.2	68.4	72.9
NRML (HOG)	83.7	74.6	71.6	80.0	77.5	80.8	72.8	74.8	70.4	74.7
BIU (LBP) [112]	85.5	76.5	69.9	74.4	76.6	84.2	79.5	76.0	77.0	79.2
BIU (HOG)	<b>86.9</b>	76.5	70.6	79.8	78.4	87.5	80.8	79.8	75.6	81.0
VGG-Face [46]	72.0	77.6	78.3	80.6	77.1	68.8	74.4	76.6	74.6	73.6
ResNet + CF [158]	78.0	<b>83.7</b>	<b>87.0</b>	<b>80.8</b>	<b>82.4</b>	<b>87.7</b>	<b>86.0</b>	<b>86.7</b>	<b>87.4</b>	<b>86.6</b>

### 5.5.5 Transfer-Learning Experiment

To demonstrate that FIW generalizes well, we fine-tune the ResNet CNN model on the entire dataset and assess the model on a smaller, non-overlapping image collection. Specifically, we achieve state-of-the-art performance using a fine-tuned CNN to encode faces of the renown KinWild datasets (see Table 5.7). For KinWild I, we get a 4% increase in performance (*i.e.*, from 78.4% to 82.4%). For KinWild II, there was a 5.6% improvement to 86.6%.

Notice the significant boost in accuracy for KinWild I for type F-D, and especially compared with M-D. Clearly, the small sample size of these types is not properly represented in KinWild, while FIW yields far less variance between scores of parent-child types. Regardless of the high score of *Bar Ilan University* (BIU) for type F-D, our fine-tuned network performs better on all other types, in average accuracy, and while providing less variation in type-specific scores. Again, this variance is caused by the small sample size, as there is less variation in score for the parent-child types of FIW.

### 5.5.6 Human assessment using FIW

We evaluated humans in kinship verification with a subset of FIW pairs. Although others conducted similar experiments [125, 139, 168], this was done with a larger sample set made up of

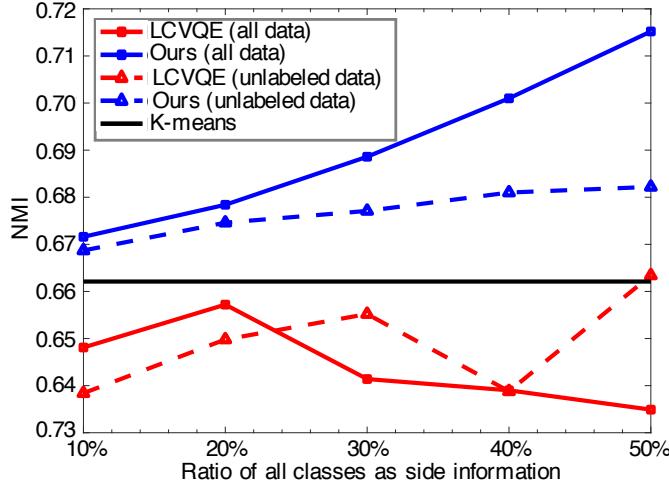


Figure 5.11: **Results for clustering families using different amounts of side information.** As clearly depicted, our method obtains the top performance. Moreover, a distinct increase in NMI for our method is shown with an increase in the amounts of side information.

more relationship types (*Case 1*). Additionally, an evaluation was done for the Boolean case only (*Case 2*). We now discuss experimental settings, results, and analyses of both human experiments.

#### 5.5.6.1 Experimental Setting

First a list of pairs from FIW with a fair data distribution was sampled (*i.e.*, different and diverse families with faces of various ages). Faces for both positive and negative pairs were from different photos. Also, we used no more than one positive and negative sample per member. We rigorously examined and, in some ways, handcrafted the list to best control the experiment (*i.e.*, replaced face images of poorer quality and famous people). Thus, efforts were spent to better ensure a fair, unbiased assessment. We also only used faces to avoid evidence besides facial appearance influencing human responses [169]. The same list of images was used for both cases: evaluating pairs per specific relationship types and for the Boolean case only.

A Google Form was used to collect responses, and the university and social media networks to recruit volunteers. Answers were anonymous, although demographic information was collected (*i.e.*, ethnicity, country of origin, and gender). Some volunteers completed both experiments. However, scores and answers were not revealed. Also, there was nearly a year between when the two experiments were conducted, with the Boolean case being a follow-up experiment to analyze

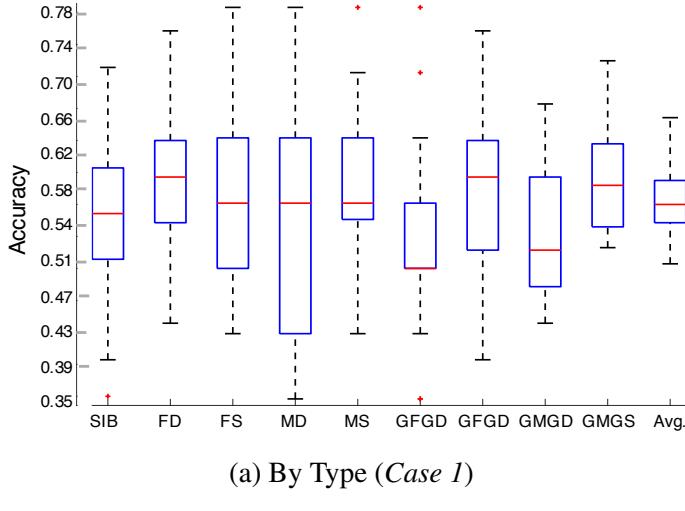


Figure 5.12: **Box plot for humans on kinship verification.** *Case 1:* Relationship type dependent evaluations. *Case 2:* Evaluations with type unspecified.

how specific relationship types influence responses. Users chose from predefined responses: *Related*, *Unrelated*, or *Skip*. Participants were asked to *Skip* if they had prior knowledge of one or both subjects, regardless of knowledge about the relationships (*i.e.*, skip any pair containing an identifiable face). Face pair-types were processed in no special order: a type-by-type basis for *Case 1*, then shuffled at random for *Case 2*. There was a total of 406 face pairs sampled from the 11 categories (see Table 5.8).

We had 75 and 110 volunteers for *Case 1* and 2, respectively. No training of any sort was provided. In both cases, the distribution of demographics was approximately 45% Caucasian, 35% Asian, 10% Hispanic/ Latino, 4% African American, and 1% Arab; 65% born in the United States, 30% from China, and 1-2% from South America, Middle East, and the Philippines; 55% males and 45% Females. No specific demographics were targeted (*i.e.*, a matter who volunteered on social media, per request of the authors, etc.). Future work could involve a greater emphasis on demographics, both in overall distribution of volunteers and intended analysis. Here, we hope to lay the framework for such a study, along with other interesting directions that assessing human ability to recognize kinship can take.

To compare human performance to benchmarks, we fine-tune SphereFace CNN on the 764 families that were not included in face pairs used for the human evaluation.

**Table 5.8: Face pair counts for human evaluation on kinship verification.** SIBS represents all siblings of the same generation, PC are parent-child, and GPGC are grandparent-grandchild.

	SIBS	PC	GPGC	Total
No. Face Pairs (per type)	50	36	28	—
No. Face Pairs (total)	150	144	112	370

### 5.5.7 Discussion

The label structure of FIW is dynamic— labels can be parsed to use the data in various ways. For instance, siblings can be split between those who share one and both parents. Even a slight change in paradigm can drastically change the study— use both parents for verification (*i.e.*, tri-subject verification [113]); use child photos only to test with for family classification. Besides, we still need to improve our visual recognition capability for kinship in current benchmarks. Then, it only seems natural to aim for fine-grained categorization of entire family trees (*i.e.*, the ultimate achievement). On a different note, generative modeling is another interesting research track to pursue (*e.g.*, given a couple and predict the offspring, or samples of their baby and predict the baby’s appearance as an adult). Even other pair types (*e.g.*, great- and great-great-grandparents, cousins, aunts, uncles, etc.). Also, the labeling framework introduced in this work could be used to add video data to the families of FIW, which can be served as a resource for template- based search and retrieval, or even consider emotional responses and facial expressions of family members.

We expect that as researchers advance this problem, FIW and its uses too will advance, and especially when considering the potential for interdisciplinary collaborations— Whether nature-based studies, generative or predictive modeling, or security-based. We hope FIW inspires new types of problems, and anticipate the list of uses to only grow when FIW is in the hands of researchers worldwide. In the end, the aim here is to attract more experts to the problem of kinship recognition.

*Families In the Wild* (FIW) is the first large-scale dataset available for visual kinship recognition. We annotated complex hierarchical relationships with only a small team in a fast and efficient manner— providing the largest labeled collection of family photos to-date. FIW was structured to support multiple tasks with its dynamic label structure. We provided several benchmarks for kinship verification and family classification. Pre-trained CNNs were used as *off-the-shelf* face encoders, which outperformed conventional methods. Results for both tasks were further improved by fine-tuning the CNN models on FIW. We measured human observers and compared their performance

Table 5.9: **Kinship verification (T1) counts.** Number of unique pairs (**P**), families (**F**), and face samples (**S**), with an increase in counts and types since [5].

	BB	SS	SIBS	FD	FS	MD	MS	GFGD	GFGS	GMGD	GMGS	Total	
Train	<b>P</b>	991	1,029	1,588	712	721	736	716	136	124	116	114	6,983
	<b>F</b>	303	304	286	401	404	399	402	81	73	71	66	2790
	<b>S</b>	39,608	27,844	35,337	30,746	46,583	29,778	46,969	2,003	2,097	1,741	1,834	264,540
Val	<b>P</b>	433	433	206	220	261	200	234	53	48	56	42	2,186
	<b>F</b>	74	57	90	134	135	124	130	32	29	36	27	868
	<b>S</b>	8,340	5,982	21,204	7,575	9,399	8,441	7,587	762	879	714	701	71,584
Test	<b>P</b>	469	469	217	202	257	230	237	40	31	36	33	2,221
	<b>F</b>	149	150	89	126	133	136	132	22	21	20	22	1,190
	<b>S</b>	3,459	2,956	967	3,019	3,273	3,184	2,660	121	96	71	84	39,743

to the machine vision algorithms, showing that CNNs surpass humans in recognizing kinship.

## 5.6 Data challenges and incentives

Challenges date back to 2011, where multi-modal data for twins was collected annually and in a highly controlled setting (*i.e.*, Twins Day [107]). Also, starting in 2014 were data challenges on unconstrained face data [120]. Then, Lu *et al.* attracted many with a IEEE International Conference on Automatic Face and Gesture Recognition (FG) challenge with KinFaceW [112]. Robinson *et al.* expanded the data challenges as part of a 2017 ACM MM Workshop using the first large-scale visual kinship recognition dataset [5], which was followed by three consecutive FG challenges - an annual effort that still occurs nowadays [7] (*i.e.*, 2018-2020, with 2020 still accessible via Codalab<sup>2</sup>). Besides, over five hundred teams partook in RFIW on Kaggle.<sup>3</sup> Recently, there have been several tutorials at top-tier conferences (*i.e.*, ACM MM18 [30], CVPR 2019<sup>4</sup>, and FG 2019<sup>5</sup>). The human evaluations were done using volunteers in a non-competitive forum.

## 5.7 Experimental

The organization of this section is as follows. First, we examine studies involving a human’s ability to recognize kinship as imagery. Thus, deeming the soft-attribute of kinship as being detectable by the eye. Next, we review kin-based task protocols - each complete with a problem statement, data splits, metrics, and baseline solutions. We then highlight commonalities in problem formulation and proposed solutions for the various tasks. Following this, we describe traditional and deep solutions. We then put this in perspective with the RFIW data challenge series - four editions (*i.e.*, 2017 [5]-2020 [7] and Kaggle Competition<sup>6</sup> held just prior to the 2020 RFIW). Finally, we discuss recent attempts to predict the appearance of family members’ faces.

### 5.7.1 Task Evaluations, Protocols, Benchmarks

We next describe each kin-based task separately: the problem statement and motivation, data splits and protocols, and benchmark experiments (*i.e.*, baselines). A brief section on the common experimental settings precedes the detailed descriptions of settings unique to the task and follow in separate subsections.

#### 5.7.1.1 Common settings

The FIW dataset provides the most extensive set of face pairs for kin-based face recognition. FIW provides the data needed to train modern-day data-driven deep models [170, 171, 29]. With over 12,000 family photos for 1,000 disjoint family trees the data contains various counts for faces, samples, members, and relationships per family. Hence, the faces of the image collection are cropped out and organized by family— faces of each family member ranges from one-to-many. FIW is split into three parts: *train*, *val*, and *test*. Specifically, 60% of the families were assigned to the *train* set; the remaining 40% was split evenly between *val* and *test*. The three sets are disjoint in family and identity. The test set remains “blind”, with automatic scoring of submissions added to the leadership board of the codalab competition. Note that the splits are consistent across tasks, so the same families makeup the “blind” set.

---

<sup>2</sup><https://competitions.codalab.org/competitions/21843>

<sup>3</sup><https://www.kaggle.com/c/Recognizing-Faces-in-the-Wild>

<sup>4</sup>[https://web.northeastern.edu/smilelab/fiw/cvpr19\\_tutorial/](https://web.northeastern.edu/smilelab/fiw/cvpr19_tutorial/)

<sup>5</sup><http://fg2019.org/visual-recognition-of-families-in-the-wild>

<sup>6</sup><https://www.kaggle.com/c/recognizing-faces-in-the-wild>

Table 5.10: **Kinship verification (T1) results.** Averaged verification accuracy scores of RFIW.

Methods	BB	SS	SIBS	FD	FS	MD	MS	GFGD	GFGS	GMGD	GMGS	Avg.
ArcFace [172] (baseline)	0.57	0.64	0.50	0.61	0.66	0.69	0.62	0.66	0.71	0.73	<b>0.68</b>	0.64
stefhoer [173]	0.66	0.65	0.76	<b>0.77</b>	0.80	0.77	<b>0.78</b>	0.70	<b>0.73</b>	0.64	0.60	0.74
ustc-nelslip [174]	0.75	0.74	0.72	0.76	0.82	0.75	0.75	<b>0.79</b>	0.69	<b>0.76</b>	0.67	0.76
DeepBlueAI [175]	0.77	0.77	0.75	0.74	0.81	0.75	0.74	0.72	<b>0.73</b>	0.67	<b>0.68</b>	0.76
vuvko [176]	<b>0.80</b>	<b>0.80</b>	<b>0.77</b>	0.75	<b>0.81</b>	<b>0.78</b>	0.74	0.78	0.69	<b>0.76</b>	0.60	<b>0.78</b>

As part of pre-processing, faces for all three sets were encoded via ArcFace CNN [172] (*i.e.*, 512 D). All pre-processing and the model weights were from the original work.<sup>7</sup> Also common is the use of cosine similarity to determine closeness of a pair of facial features  $p_1$  and  $p_2$  [177]. This is defined as

$$CS(p_1, p_2) = \frac{\mathbf{p}_1 \cdot \mathbf{p}_2}{\|\mathbf{p}_1\| \cdot \|\mathbf{p}_2\|}.$$

Scores were then either compared to threshold  $\gamma$  (*i.e.*,  $CS(p_1, p_2) > \gamma$  infers *KIN*; else, *NON-KIN*) or sorted (*i.e.*, ranked list). This concludes the common experimental settings.

Teams were allowed up to six final submissions per task. Submissions were accompanied by a brief (text) description of the system used to generate results.

### 5.7.1.2 Kinship verification (T1)

Kinship verification aims to determine whether a face pair are blood relatives. This classical Boolean problem has two possible outcomes, *KIN* or *NON-KIN*. Hence, this is the *one-to-one* view of kin-based problems. The classical problem can be further extended by considering the type of kin relation between a pair of faces, rather than treating all kin relations equally [30].

Prior research mainly considered parent-child kinship types, *i.e.*, FD, FS, MD, MS. Less attention has been given to sibling pairs, *i.e.*, SS, BB, and SIBS. Research findings in psychology and computer vision found that different relationship types share different familial features [105]. Hence, each relationship type can be modeled and evaluated independently. Thus, additional kinship types would further both our understanding and capabilities of automatic kinship recognition. With FIW, the number of facial pairs accessible for kinship verification has dramatically increased. Additionally, benchmarks now include grandparent-grandchildren types, *i.e.*, GFGD, GFGS, GMGD, GMGS.

#### Data splits.

---

<sup>7</sup><https://github.com/ZhaoJ9014/face.evoLVE.PyTorch>

## CHAPTER 5. FAMILIES IN THE WILD (FIW)

The two datasets used (*i.e.*, KinFaceW and FIW) follow the same settings. Both provide a list of pairs labeled as KIN and NON-KIN. The differences are in the number of pair types and overall size of splits. Data specifications are in Table 5.9.

KinFaceW provides two sets (*i.e.*, KinFaceW I & II) and the four parent-child pair types. FIW spans eleven different relationship types - the types used in 2020 RFIW (Table 5.2). The test set is made up of an equal number of positive and negative pairs and with no family (and, hence, no identity) overlap between sets.

### Settings and metrics.

Verification accuracy is used to evaluate. Specifically,

$$\text{Accuracy}_j = \frac{\# \text{ correct predictions for } j\text{-th type}}{\text{Total } \# \text{ of pairs for } j\text{-th type}},$$

where  $j \in \{4 \text{ relationship types and } \emptyset \text{ for KinFaceW and } 11 \text{ relationship types and } \emptyset \text{ for FIW}\}$  (listed in Fig. 5.13). Then, the overall accuracy is calculated as a weighted sum (*i.e.*, weight by the pair count to determine the average accuracy).

**Baseline and results.** The threshold was determined by the value that maximizes the accuracy on the *held-out* data in all cases. The results on KinFaceW I and II are show in Table 5.7. The results for FIW are in Table 5.10, with sample pairs that either 100% or 20% of all teams got correct are shown in Fig. 5.14.

#### 5.7.1.3 Tri-Subject verification (T2)

Tri-Subject Verification (T2) focuses on a different view of kinship verification– the goal is to decide if a child is related to a pair of parents. First introduced in [113], it makes a more realistic assumption, as having knowledge of one parent often means the other potential parent(s) can be easily inferred.

Triplet pairs consist of Father (F) / Mother (M) - Child (C) (FMC) pairs, where the child C could be either a Son (S) or a Daughter (D) (*i.e.*, triplet pairs are FMS and FMD).

**Data splits.** Following the procedure in [113], we create positive (have kin relation) triplets by matching each husband-wife spouse pair with their biological children, and negative (no kin relation) triplets by shuffling the positive triplets until every spouse pair is matched with a child which is not theirs (Table 5.11).

The number of potential negative samples far exceeds the number of potential positive examples– We post-process the positive triplets before generating negatives to ensure balance among

individuals, families, and spouse pairs, since a naive data selection procedure which weights every face sample similarly would result in some individuals and families being severely over-represented due to an abundance of face samples for some identities and families. The post-processing is done by limiting the number of samples of any triplet  $(F, M, C)$ , where  $F$ ,  $M$ , and  $C$  are identities of a father, mother, and child to 5, then limiting the appearance of each  $(F, M)$  spouse-pair to 15, and then finally limiting the number of triplet samples from each family to 30. The *test* set has an equal number of positive and negative pairs. Lastly, note that there is no family or subject identity overlapping between any of the sets.

**Settings and metrics.** Per convention in face verification, we offer 3 modes (*i.e.*, the same as in task 1 listed in Section 5.7.1.2). Again, the metric used is verification accuracy, which is first calculated per triplet-pair type (*i.e.*, FMD and FMS). Then, the weighted sum (*i.e.*, average accuracy) determines the leader-board.

**Baseline and results.** Baseline results are shown in Table 5.12, with samples of easier and more challenging samples for both *KIN* and *NON-KIN* triplets in Fig. 5.15 and 5.16. A score was assigned to triplet  $(F_i, M_i, C_i)$  in the validation and *test* sets using the formula

$$\text{Score}_i = \text{avg}(\cos(F_i, C_i), \cos(M_i, C_i)),$$

where  $F_i$ ,  $M_i$  and  $C_i$  are the feature vectors of the father, mother, and child images respectively from the  $i$ -th triplet. Scores were compared to a threshold  $\gamma$  to infer a label (*i.e.*, predict *KIN* if the score was above the threshold; else, *NON-KIN*). The threshold was determined experimentally on the *val* set and used for *test*.

#### 5.7.1.4 Search and retrieval (T3)

A newer task, search and retrieval (T3), is posed as a *many-to-many*, *i.e.*, one-to-many samples per subject. Thus, we imitate template-based evaluations on the probe side, but with the gallery now labeled by family. Furthermore, the goal is to find relatives of search subjects (*i.e.*, *probes*) in a search pool (*i.e.*, *gallery*).

The protocol of T3 could be used to find parents and other relatives of unknown, missing children. The gallery contains 31,787 facial images from 190 families (Fig. 5.17): inputs are subject labels (*i.e.*, probes), and outputs are ranked lists of all faces in the gallery. The number of relatives varies for each subject, ranging anywhere from 0 to 20+. Furthermore, probes have one-to-many

samples— the means of fusing samples of probes is an open research question. This *many-to-many* task is currently set up in closed form (*i.e.*, every probe has a relative(s) in the gallery).

**Data splits.** This task will be composed of search subjects (*i.e.*, *probes*) from different families. *Probes* are supported by several samples of query subject, text description of family (*e.g.*, ethnicity, some relationships between selected members, etc.), and list of relatives present in the *gallery*. The *test* set will only consist of sets of images for the probes. Diversity in terms of ethnicity is ensured for both sets. Again, three disjoint sets were split (Table 5.13).

**Settings and metrics.** Each subject (*i.e.*, probe) is searched independently— 190 subjects with one-to-many faces. Hence, 190 families make up the *test* set. Following template conventions of other *many-to-many* face evaluations, face images of unique subjects are separated by identity, with a gallery containing various number of relatives with a variable number of faces each [178].

Mean average precision (MAP) was the underlying metric used for comparisons. Mathematically speaking, scores for each of the  $N$  missing children are calculated as follows:

$$AP(f) = \frac{1}{P_F} \sum_{tp=1}^{P_F} Prec(tp) = \frac{1}{P_F} \sum_{tp=1}^{P_F} \frac{tp}{rank(tp)},$$

where average precision (AP) is a function of family  $f$  with a total of  $P_F$  true-positive rate (TPR). MAP is then

$$MAP = \frac{1}{N} \sum_{f=1}^N AP(f).$$

**Baseline and results.** Submissions consisted of a matrix with a row per *probe* listing the indices of all subjects in the *test* gallery as a ranked list. Results are listed in Table 5.14 with sample inputs and predictions shown in Fig. 5.18.

## 5.8 Methodologies

Many formulated kinship recognition problems in the visual domain as multi-view, multi-task, and multi-modal, which is typically to increase the amount of information obtainable, even when the final target is among other targets during training (*i.e.*, auxiliary tasks that complement the knowledge obtained from recognition, alone). For instance, the Deep Kinship Matching and Recognition (DKMR) was proposed as a jointly-trained model on top of a graph optimization algorithm [179]. Clearly, deep learning has overcome the traditional metric-learning approaches from about 2017 [171, 180, 181, 182, 183] and still today [184, 185, 186, 187, 188]. We will first

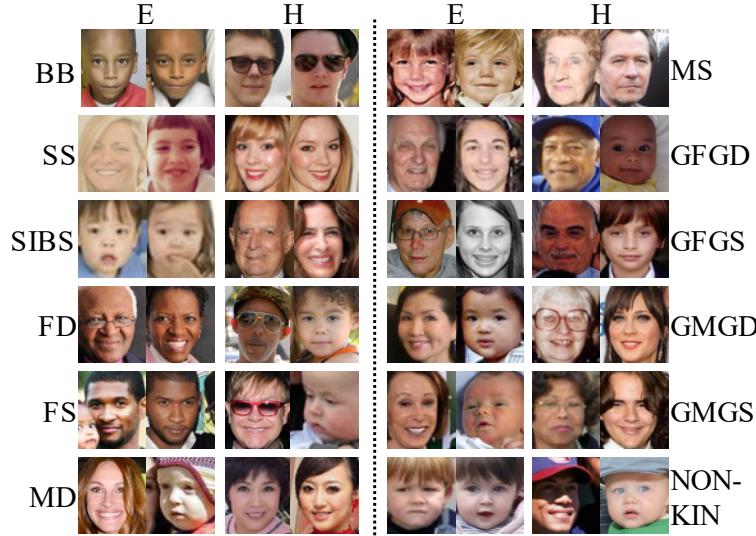


Figure 5.13: **Kinship verification (T1) sample pairs.** Sample pairs with similarity scores near the threshold (*i.e.*, hard (H) samples), along with highly confident predictions (*i.e.*, easy (E) samples) in verification task.

review the traditional methods, and then deep learning, for discriminating problems; an overview of the kin-based generative modeling is given at the end of the section.

Table 5.11: **Tri-subject verification (T2) counts.** No. pairs (P), families (F), face samples (S).

		FM-S	FM-D	Total
train	<b>P</b>	662	639	1,331
	<b>F</b>	375	364	739
	<b>S</b>	8,575	8,588	17,163
val	<b>P</b>	202	177	379
	<b>F</b>	116	117	233
	<b>S</b>	2,859	2,493	5,352
test	<b>P</b>	205	178	383
	<b>F</b>	116	114	230
	<b>S</b>	2,805	2,400	5,205

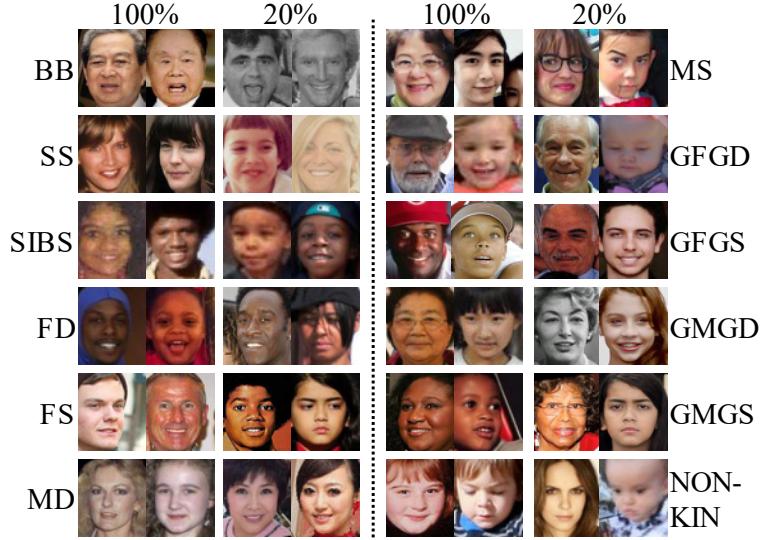


Figure 5.14: **Qualitative analysis of T1.** Samples of each relationship type that all of the teams either got correct (100%) or mostly not (20%) for the eleven pair types of FIW and NON-KIN.

### 5.8.1 Traditional approaches

The main focus of the survey is on large data resources, along with the modern-day complex, data-driven modeling (*i.e.*, deep learning). However, such respective work makes up the latter half of the decade. Feature and metric-learning dominated the first half of this past decade in research of visual kinship recognition— before the release of FIW. For completeness, we will introduce several methods that predate the deep methods on FIW.

**Handcrafted features.** Fang *et al.* proposed using features such as geometric differences between face parts, color features, and handcrafted features that were the basis for the metrics to be learned in the years to come [34]. Furthermore, and as mentioned, many of the smaller datasets are limited in diversity (*i.e.*, all similar demographics) and with pairs from the same photos, from which some proposed color-based features [189]. Still, papers that hone-in on the smaller data employ more classical approaches, such as representation learning via binary trees [190].

**Metric learning.** Metric learning methods are popular solutions in kin-based vision problems. The general idea is to optimize a metric between classes. In kinship verification, the classes are KIN and NON-KIN (*i.e.*, true match and imposter, respectively). Lu *et al.* proposed NRML for kinship verification which aims for a contractive deep belief net (fcDBN) made by stacking fcRBMs to learn weights in a greedy, layer-by-layer fashion using both local and global features [132].

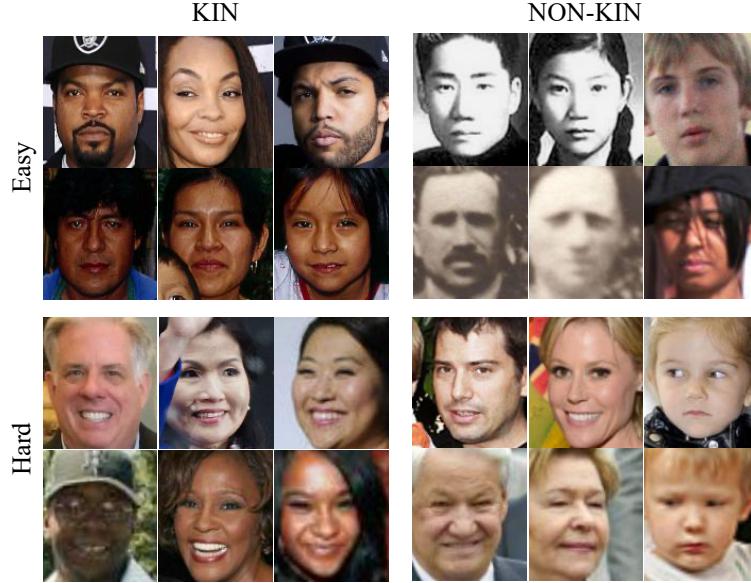


Figure 5.15: **Triplets with extreme scores (i.e., correct and incorrect).** Each show FMS (top rows) and FMD (bottom) for tri-subject (T2).

Wu *et al.* combined color and texture features for kinship verification with extreme learning machines (ELM) for robustness on small data [131]. Mahpod *et al.* proposed a hybrid asymmetric distance learning (MHDL) scheme, combining symmetric and asymmetric multiview distances [191]. Most recently, Hu *et al.* proposed treated-different features as multiple views via a multi-view geometric mean metric learning (MvGMML) [133].

For more details on the traditional methods see [118].

Table 5.12: **Verification scores.** Results for tri-subject (*i.e.*, T2).

	FMS	FMD	Avg.
Spherefase [172] (baseline)	0.68	0.68	0.68
stefhoer [173]	0.74	0.72	0.73
DeepBlueAI [175]	0.77	0.76	0.77
ustc-nelslip [174]	<b>0.80</b>	<b>0.78</b>	<b>0.79</b>

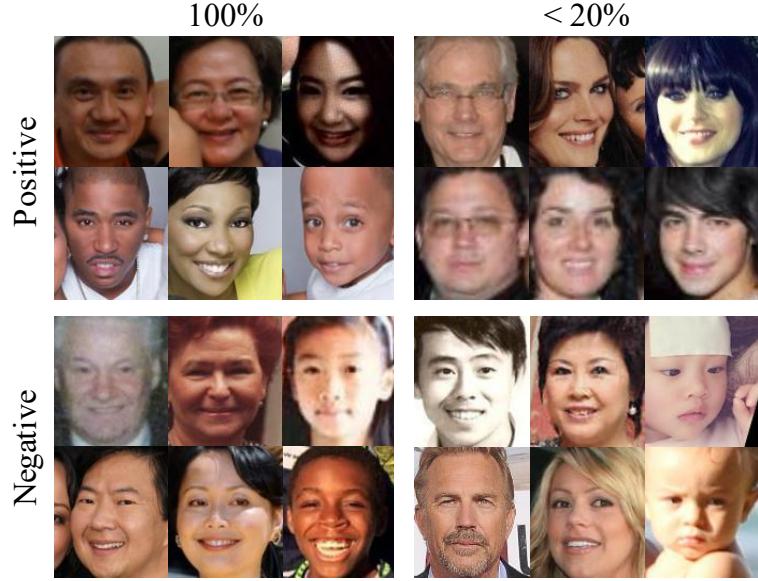


Figure 5.16: **Sample of T2.** Samples that all teams got correct (left) and mostly incorrect (right) for FMS (top rows) and FMD (bottom).

### 5.8.2 Deep learning approaches

The 2012 AlexNet [43] sparked the deep learning era. As done in many problems, deep learning grew more popular with the big data provided with FIW. Still on KinFaceW, we first review the deep metric done using small amounts of training data, and then discuss the data-driven work done using FIW.

There are many commonalities between the different solutions proposed as part of the RFIW challenge. Typically, a ResNet-based [16] backbone; if not, then together with FaceNet [142]. Nonetheless, the story as seen in the timeline is split in half (*i.e.*, with the latter half dominated by modern-day deep learning approaches) and quite significantly, metrics learned on top of hand-crafted features dominated the charts as SOTA for many years [192]. As recent as 2017, metric-learning was a go-to approach for kin-based problems, whether a single metric or multiple (*e.g.*, large-margin multi-metric learning (LM<sup>3</sup>L) [193]). Even so, geometric and distant features in pixel space (*e.g.*, key point coordinates on neutral face [194])—directly related to insufficient data for modern-day data-driven machinery (*i.e.*, deep learning).

Provided a deep CNN trained to classify face identity, the encodings produced encapsulated much information of the subject. However, instead of looking for absolute closeness in embedding space as the ideal case for a set of samples of a single class (*i.e.*, identity), in kin-based tasks we hope

Table 5.13: **Tri-subject (T2) counts.** Individuals **I**, families **F**, face samples **S**.

		Probe	Gallery	Total
train	<b>I</b>	–	3,021	3,021
	<b>F</b>	–	571	571
	<b>S</b>	–	15,845	15,845
val	<b>I</b>	192	802	994
	<b>F</b>	192	192	192
	<b>S</b>	1,086	4,030	5,116
test	<b>I</b>	190	783	9d73
	<b>F</b>	190	190	190
	<b>S</b>	1,487	31,787	33,274

to detect when similarities between a pair (or group) of faces (*i.e.*, encoded) reflect that of the various relationship-types. For this, many tend to fine-tune models initially trained on a larger FR-based database, such as VGG-Face [142], VGG2 [195], and MSCeleb [196]. We next speak on various flavors of deep learning.

**Pre-trained CNNs.** Besides that most solutions involve the renowned Siamese training model, many of which still incorporate a cosine loss as in the seminal work done at Bell Lab’s mid-90s [197], *i.e.*, multiple inputs to networks with shared weights for which metric is learned on top (Fig 4.5). In the simplest form, Siamese-based CNN models map two or more samples by a single CNN to a real-number vector space  $\mathbb{R}^d$  (*i.e.*, a function  $f(\cdot)$ ) to encode an image (*i.e.*, facial [encoding, embedding, feature] of size  $d$ , especially in the context of facial representation, all refer to the  $f(x_i) = z_i \in \mathbb{R}^d$ ). Generally, and in most methods proposed in RFIW, the shared model is pre-trained data for another, yet similar task (*i.e.*, facial recognition). With that, the CNN that now serves as an encoder, maps  $k$  samples to its  $d$ -dimensional space learned to discriminate between faces. With the Siamese frozen—whether entire network, with a couple of layers on top set with a small learning rate, or popped off by adding a path that splits off prior to later rejoin or just remove entirely—the goal then is to learn a metric optimal for recognizing family members by face cues. Clearly, there are several design choices—with simple solutions in those with an *off-the-shelf* CNN with no additional training (*i.e.*, trained for FR, so naively assuming that the best way to detect kinship is to detect faces that look like the source). However simple, and with many cases a fair assumption, the naive approach

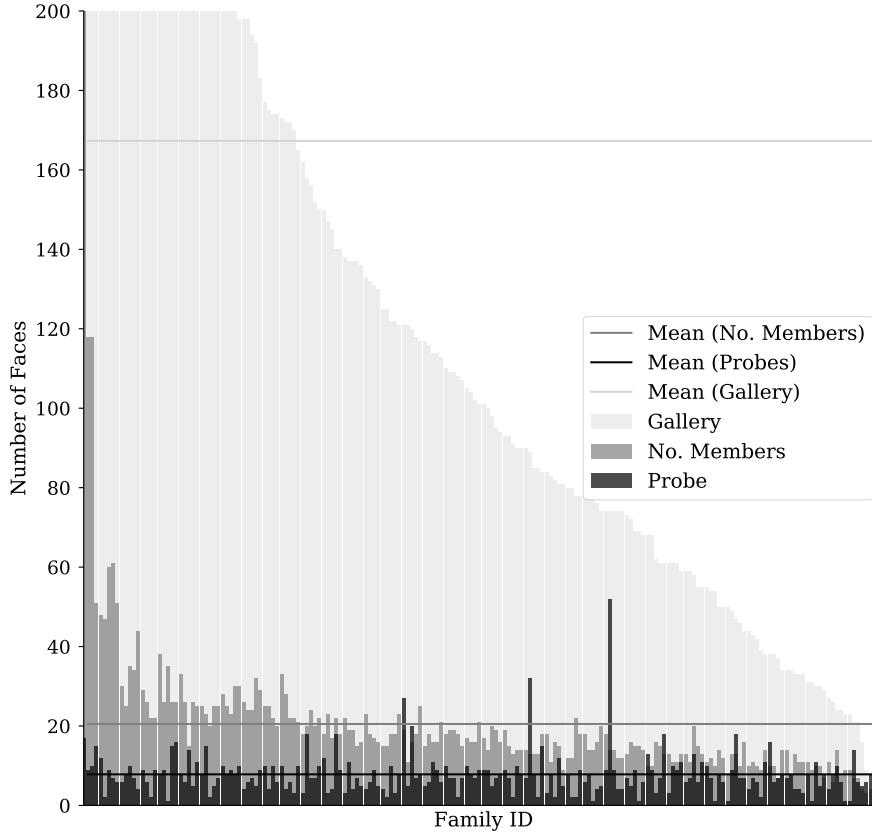


Figure 5.17: **Plot of face counts per family in test set of T3.** The probes have about 8 faces on average, while the number of family members in the gallery nears 20 on average, with a total average of 170 faces.

outperformed previous SOTA methods prior to FIW providing the number of data samples needed to suffice the capacity of most deep learning approaches. In light of this, the CNN then serves as the method for feature extraction— claiming to provide the best face representations for the task. As previously described of the wave of metric-based and subspace-modeling methods, we can then further refine the output of the feature extractor by extending the composition function by adding and training mappings in the embedding space and while again, often with the weights of the pre-trained CNN  $f$  held static. From this, kin-based tasks can be targeted by learning filters, mappings, and even metrics from the embedding space on up (*i.e.*, build up from the embedding space from where face embeddings are compared in some fashion).

**Deep metric learning.** Lu *et al.* proposed to learn a distance metric for  $K$  feature types via  $K$  MLPs - learn to project each feature using the optimal thresholds determined independently [134]. This

Table 5.14: **T3 results.** Performance ratings for SOTA methods.

Methods	mAP	Rank@5
Baseline (Sphereface) [172]	0.02	0.10
DeepBlueAI [175]	0.06	0.32
HCMUS notweeb [202]	0.07	0.28
ustc-nelslip [203]	0.08	0.38
vuvko [176]	<b>0.18</b>	<b>0.60</b>

method, which was called discriminative deep metric learning (DMML), proved effective on the KinFaceW settings of minimal training data (Table 5.7).

**Fine-tuning.** There is an abundant of public FR data (*e.g.*, LFW, VGG, MSCeleb [196]) with some labeled by *soft attributes* (*e.g.*, age [198], gender [199], attribute, and diverse demographics [25, 200]). With this, and provided the known concept of deep learning tending to learn transferable features [201], the use of fine-tuning pre-trained has been done by many. For instance, a SphereFace loss, which is a multi-class loss, is first used to train a large CNN to do facial recognition on identities of an auxiliary dataset, and then having the layers near the top fine-tuned to recognize the families of the FIW training set via

$$\mathcal{L}_{\text{family}}(\theta) = -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp^{W_{y_i}^T x_i + b_{y_i}}}{\sum_{j=1}^N \exp^{W_{y_j}^T x_i + b_j}}, \quad (5.7)$$

where  $B$  is the batch size,  $N$  is the number of families,  $x_i$  is the face encoding from family  $y_i$ ,  $W$  is the weight matrix (*i.e.*,  $W_j$  denotes the  $j^{\text{th}}$  column) and  $b$  is the bias term. In the end, verifying kinship between a face pair can be done using the model to encode the faces and cosine distance to measure their closeness. If family labels are unavailable, which is another setting of the verification task, approaches tend to use Siamese concepts on top of the pre-trained CNN (Fig. 4.5). Specifically, sharing weights for two or more samples, and penalizing based on the closeness between a set of samples upon being encoded by the network, has shown to be an effective means of staging a network for the verification task. In return, Siamese; furthermore, the relationship between the pairs with respect to labels at training differences is in preprocessing, method of fusion (*e.g.*, *early* versus *late*).

In [203], Track I and III completed in succession, such that a wider sweep of CNN backbones, loss functions, and fusion methods were assessed in Track 1, to both gain deeper understanding to make decisions pertaining to Track III. Mainly, ResNet50 and SENet50 were

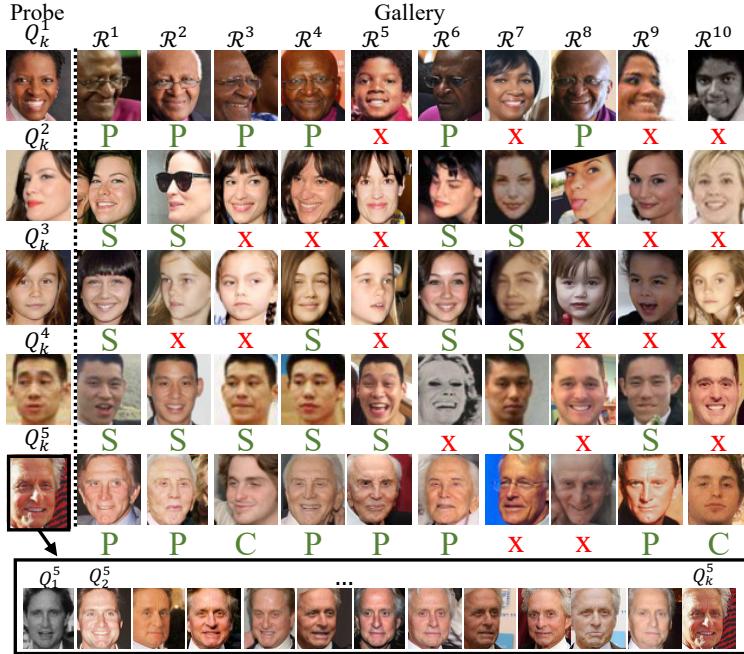


Figure 5.18: **T3 sample results (Rank 10).** Each query (row) has one or more faces, for the probe returns and ranks all samples in the gallery - here we show top 10. FP are labeled by **X**, while true matches list the relationship type in green: **P** for parent; **C** for child; **S** for sibling.

evaluated separately, each with additional fully-connected layers with two losses on top, Binary Cross Entropy (BCE) and Focal loss. BCE, a widely used loss that does as its name implies: uses the measure of entropy of a distribution, say  $q(y)$  for  $c \in 1, \dots, C$  classes as  $H(q) = \sum_{c=1}^C q(y_c) * \log(q(y_c))$ . Since we have no knowledge of the true distribution, we aim to match samples of the *true* distribution  $p(y)$ . Hence, cross-entropy is entropy between  $p(y)$  and  $q(y)$ .

Yu *et al.* found that BCE loss outperformed Focal Loss for all fusion schemes and settings in Track I [203]. Intuitively, this makes sense as Track I, a Boolean task, has an equal number of positive and negative pairs – imbalanced data motivated Focal Loss, which is not an issue for verification. Then, transferring over the model, loss, and fusion settings that worked best for Track I to Track III and used as is. The difference is in the ranking scheme (*i.e.*, provided multiple faces per query, the average of all faces and each gallery sample determined the score at the subject-level).

**Deep representation.** Training a set of CNNs, each targeting specific regions (or parts) of the face, was proposed as a solution for KinFaceW [204]. Then, Heterogeneous Similarity Learning (HSL) tackled various tasks of kinship recognition via multi-view learning, with the different views set as different relationship types dubbed multi-view SML (MSML). [180]. Similarity, Support Vector

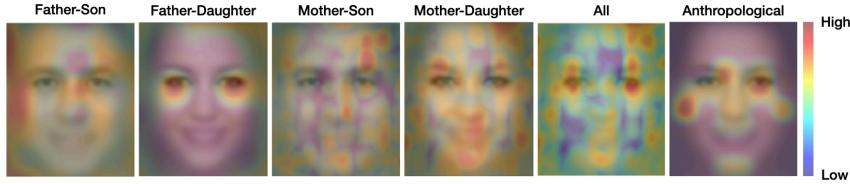


Figure 5.19: **Activations from mapping image-to-latent space (from [8]).** The salience mapped from the activation response and superimposed on the average face. Family101 dataset was used for this experiment [9]. The end result depicted here were dubbed the *genetic features* from latent space of a trained Gated AE.

Data Description (SVDD) was proposed as a Support Vector Data Description-based metric learning (SML) loss function, allowing detailed information to be extracted as geometric and appearance-based features for kinship verification [205]. Duan *et al.* proposed a coarse-to-fine scheme for which CNNs at different levels (*i.e.*, layers) were transferred from being trained using a FR dataset and then fine-tuned for kinship using a loss function based on NRML [206]. In fact, many recent works leveraged existing FR methodologies (*e.g.*, CNN trained to classify faces) as a prior, then fine-tune using the kin-based image data as the source in a transfer-learning regime [207].

Several lines of research specifically focused on the one-to-one kinship verification problem by learning a face encoder robust in detecting kinship relationship via AE (*e.g.*, DAE [11, 208, 209]), *i.e.*, deep representation learning methods [210]. Dehghan *et al.* was amongst the first, proposing to train a Gated AE to encode faces as *genetic features*, and weighting according to the salience for the respective relationship type [8]. Fig. 5.19 depicts the salience, with high being most similar regions and low dissimilar. Besides still-faces, deep learning approaches were also proposed for recognizing kinship pairs using facial cues in video data [211]. A sequence recurrent NN was trained for kinship verification in videos using a novel attention mechanism [212]. With videos, there comes more bits of information; however, the range of bits (*i.e.*, the underlying variation of the data) should be optimized to maximize the information gain. In other words, video data introduces another space for fusion in the choosing of the best frame(s) to describe and represent [213]. Graphical neural network (GNN) with a metric learned on top proved to be one of the most effective deep learning models employed for kin-based vision problems [114]. Not just on the large-scale FIW data, but a graph-based kinship reasoning (GKR) network proved effective on KinFaceW [214] (Table 5.7).

**Approaches to data challenges.** RFIW serves as a platform for experts to publish and junior scholars to get started. The first edition of RFIW was in 2017 [5] - a data challenge workshop in conjunction

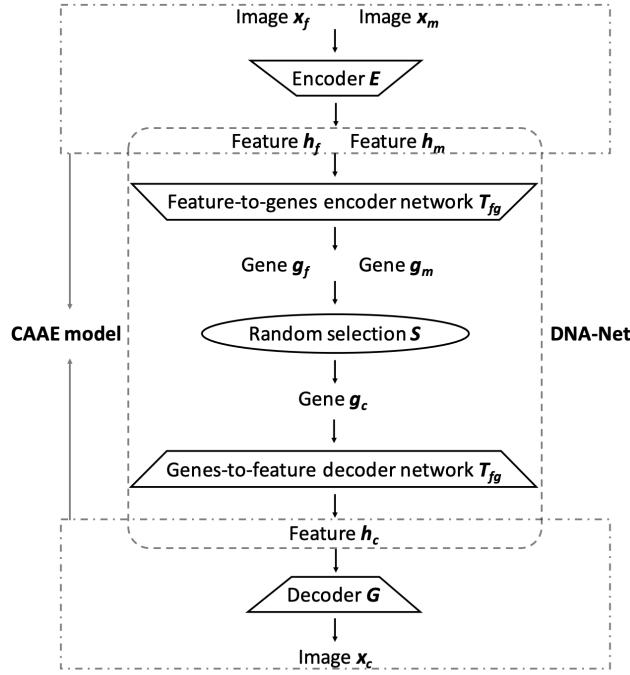


Figure 5.20: **Model to synthesize children faces from a parent-pair (visualizations from [10]).** Notice that the output of encoder  $E$  is the concatenation of features from prospective parents, the father  $h_f$  and mother  $h_m$  joined by  $\oplus$  such that the two embeddings encoded by the Siamese network are fused (*i.e.*,  $2 * \mathbb{R}^d \rightarrow \mathbb{R}^{2d}$ ) before passed as input to the CAEE model.

with the *ACM Conference on Multimedia*. Ever since, RFIW has been held annually (*i.e.*, 2018-2020 held in conjunction with FG as a data challenge), with each year building on the prior. Let us review series highlights over the years, and then focus on the top teams of the 2020 edition.

From the start, solutions for RFIW typically involved CNNs pre-trained for FR. For the top performing submission of the 2017 RFIW, Yong *et al.* used an ensemble of deep CNNs with data augmentation and mining techniques called KinNet [171]. Specifically, the authors proposed to train four ResNet models (*i.e.*, 80, 101, 152, and 269 layers) for FR to then fine-tune for kinship verification via a triplet loss targeting intra-family relationships. KinNet used two tricks during training: (1) augmentation using imaging processing techniques (*e.g.*, gamma correction, down/up sampling of pixels, blurring) and hard-negative mining for selecting triplets. In the end, KinNet scored an impressive average of 74.9%. It is important to note that the data has changed since this first edition of RFIW (*e.g.*, *grandparent-grandchild* types were not included). Thus, a comparison with the proceeding years would be unfair.

CHAPTER 5. FAMILIES IN THE WILD (FIW)



(a) Random synthesis.



(b) Age (*i.e.*, 10, 20, 30 years old from row 4-6, respectfully).



(c) Gender (*i.e.*, male-to-female from row 4-5, respectfully).

Figure 5.21: **Synthesized results from [10].** Columns correspond to families, with fathers on first row, mothers on second, and real and generated children on third row and bottom, respectively (a). See subcaption for specifics on age (b) and gender (c).

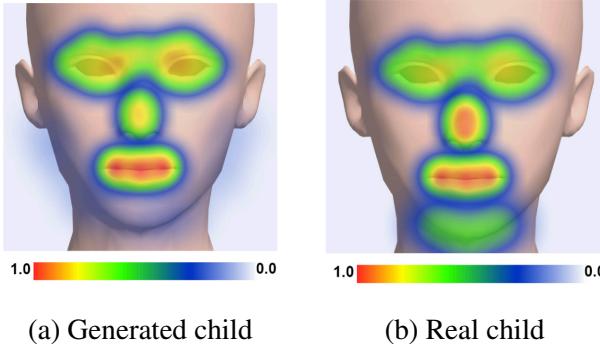


Figure 5.22: **Salience map per key-points [10]**. Best viewed in color.

In 2018, Dahan *et al.* got the top performance 68.2% [215]. Specifically, the authors trained a VGG-Face model with the novel *local features conv-layer* that fused the Siamese inputs by summing the features. In other words, conventional conv-layers share weights in image space, whereas these authors proposed learning local weights to produce pair and location specific features.

In 2019, Laiadi *et al.* extended XQDA-to-TXQDA to operate on multilinear data in a low dimensional and discriminative tensor subspace. TXQDA uses multilinear projections of tensors to a space with greater separation between data classes, is enhanced in a way and helps lightening the small sample size problem (*i.e.*, results for both KinFaceW and FIW) [216]. Nandy *et al.* followed a Siamese learning approach [217], which we will next learn was the most common in 2020.

The 2020 edition saw a great increase in interest and participation. Several methods were used as solutions for two or more tasks. Shadrikov *et al.* treated the different pair types as multiple tasks, training a local expert for each on top of a ResNet50, simultaneously. The authors used T1 data as validation, but then deployed on T2 and T3 as well. Another multi-task approach applied different fusion techniques in deep feature space [174, 203]. Zhipeng *et al.* used two pre-trained CNNs, fused the two face encodings by different types of arithmetic, and generated solutions for all three tasks [174, 203]. In [175], the distance between faces was then determined euclidean distance (instead of the typical cosine similarity); also, SENet [218] was used as the backbone showing a modest boost over ResNet50 on the validation, but dropping on the test. Like in [6], except the authors now used Arcface instead of Sphereface, Höörmann *et al.* fine-tuned a CNN using families as the classes[175] - ultimately placing second in verification (*i.e.*, T1) and tri-subject (*i.e.*, T2). Yu *et al.* put emphasis on the dependence of family identification accuracy for cross-gender versus same-gender pairs types [173]. These researchers constructed a Kinship *comparator* module that consisted of eleven “local expert networks” connected in series– eleven networks corresponding

to the eleven relationship types of T1. In the end, Stefhoer registered the highest score in the subcategories of father-daughter and mother-son in T1. Yu *et al.* also used a Siamese network, *i.e.*, encoded features from face images via a CNN with shared weights [174]. ResNet50 or SENet50 was used as the backbone, both pre-trained on VGGFace2 [195]. Team ustc-nelslip also employed two loss functions, binary cross-entropy and focal loss, and fused the features using two algebraic formulae leading to  $2 \times 2 \times 2 = 8$  independent “models.” A unique feature was the construction of a “jury system” to combine outputs of different models to improve accuracy. With [176] the top-scorer in T2. Nguyen *et al.* competed in T1 and T3[202]– the authors use a Siamese CNN with FaceNet (Inception-ResNet-v1 trained with triplet loss) and VGG-Face (Resnet-50) pre-trained. The authors also implement ArcFace [172] - a family of loss functions based on the geodesic distance between feature vectors which aim to discriminate the latent representation of deep NNs. Samples that were unanimously classified correct or most incorrect are in Fig. 5.14, 5.16, and 5.18, along with average performance ratings in Table 5.10, 5.12, and 5.14 for T1, T2, and T3, respectfully.

### 5.8.2.1 In summary

Of all the proposed methods, there is a common factor: the larger the age gap the higher percentage of FP during evaluation. As mentioned, this was addressed early on with UB Face [105] and, although fundamental to the analysis of results over the years, proposed models tended to acknowledge this as a challenge, but with no added mechanism to make robust to age-variations between *parent-children*. That is, until Wang *et al.* proposed using GAN technology to synthesize younger versions of an input face. Specifically, and a clearly effective data augmentation approach, the authors trained generators for both genders to account for this while training a deep CNN with a maximum margin loss to do boolean classification (*i.e.*, *KIN / NON-KIN*). As formalized in their work, domain  $A$ , for *aged*, was the source and domain  $Y_m$ , for *young*, was the target. Provided paired data, the *parent* aimed to transform  $x_i \in A \rightarrow x_j \in Y$  with data distribution  $x \sim p_A \rightarrow x \sim p_Y$ . Having noticed that FIW, which most closely matches real-world data, does not necessarily have parents at older ages (*i.e.*, *aged*), thus, the inputs could very well be parents as juveniles, or even during infancy. To mitigate the problem, with no age-labels provided in FIW, focus was directed to constrain the output such to influence younger aged faces less so, than if faced with an elderly parent.

### 5.8.3 Generative modeling approaches

The dynamics of the offspring synthesis problem has a great distinction from traditional one-to-one mapping - two parents with directional relationships are input as prior knowledge to predict the appearance of their child. Such a two-to-one problem raises the question on how to best fuse knowledge from a pair of faces. Let us now even consider information for various family members - the fusion then should consider directed relationships inherent to family trees. Current face synthesizers conditioned on kinship assumes knowledge of one [219] or both [10, 220] parents.

Ozkan *et al.* proposed KIN-GAN to synthesize a child's face from a sample of a single parent [219]. The problem is inherently difficult, for the variation embedded in many complex factors nearly changes from one sample-to-the-next. Nonetheless, trying to solve the problem with just one parent is insensible— it takes two to tango in nature and, thus, such a formulation is out of scope before the problem is even started. Noticing this, [220] proposed means of modeling as a two-to-one mapping. Similarly, Gao *et al.* aimed to mimic the nature of reproduction with a model dubbed DNA-Net [10]. DNA-Net fuses latent representation of a parent pair at the feature-level, which is used as input to CAEE model trained on top (Fig 5.20). The parents' signals are fused at the output of encoder E by concatenation of their features, and are then fed to the CAEE model to produce a single feature representing the face encoding of the child. Finally, the child's encoding is decoded by G to the predicted facial image.

Note that DNA-Net was dubbed by the authors in the effective work proposed; however fair when speaking in general terms (*i.e.*, infrequent situation in research), we suddenly see naming schemes such as this, *genetic features*, among few others is too strong. Nonetheless, there is a clear analogy, so for the sake of story-telling and system depiction, Gao *et al.* dubbed this as a single face is synthesized from face pair. The choice in CAEE made it so the generator could synthesize children as a function of age and sex (Fig. 5.21). Note that treating sex as a continuous spectrum, opposed to discrete labels, is both appropriate and more precise (*i.e.*, provided an extreme pair, one female and the other male, there exists many cases in between, which is, in fact, where most of society falls [221]). As a part of the work to support DNA-Net, the authors compared salience in detecting kinship of type *parent-child* at specific facial features (*i.e.*, eyes, nose, mouth, and chin). Hu invariant moments were used as the shapes of the four facial parts [222], from which the accumulative cosine distances yielded *heritability maps* (Fig. 5.22).

# Chapter 6

## FIW-MM

### 6.1 Overview

So far, we have covered vast works in automated kinship recognition that assumes a genetic relatedness between individuals detectable by facial cues - a state in technology unimaginable just over a decade ago. Much of the progress in the difficult tasks of kin-based recognition was by the availability of labeled family data with sufficient counts and concurrent advances in face recognition [48, 53]—proposed systems inherently gain if it based on a FR model that experienced a gain itself via progress in conventional FR. In other words, FR (*i.e.*, determine if face pair are of the same identity) and visual kinship recognition (*i.e.*, determine whether a pair of faces are of the same bloodline) both target facial cues to determine whether or not a face is a match to a gallery (*i.e.*, test sample(s)). Conventional FR is the more general, simpler by definition and protocols, and with a higher relevance to vast use-cases of the two. FR also has a data need that is more readily accessible. So, it is absolutely to no surprise that FR tends to be ahead of kinship recognition technology, which results in there being lots of research findings in FR that can be transferred or referenced when devising hypotheses in kinship problems.

The seminal work in visual kinship recognition introduced the first image dataset [34]. Thereafter, larger and more difficult datasets were released, such as FIW [30] and TSKIN [113]. In response, vision researchers developed methods and models to match the rising level of difficulty in kinship datasets [7, 18].

Along with conventional FR and the different sub-problems that model visual knowledge from facial cues, speaker-based problems have recently grown popular based on audiovisual data (*e.g.*, speaker separation [223], speaker identification [224, 225], cross-modal audio-to-visual or vice-

## CHAPTER 6. FIW-MM

versa [226], emotion recognition in multimedia (MM) [227, 228], and several others [229, 230]). The sudden surge of attention to audio-visual data has brought together experts who specialize in biometric signals to share thoughts, combine knowledge, and propose solutions that best fuse multi-domain knowledge for optimal decision making [231, 232]. This work shows that the addition of MM for kin-based recognition can improve the current SOTA.

Studies in speech recognition found that identical twins were hard-negatives and confused in classification tasks consistently [233]. Although results were strongly in favor of the hypothesis that twins sound alike, the experiments were done on a small sample set of 8 pairs. Interestingly, [233] showed similar results for 2 of the 8 pairs the authors dubbed ‘separated twins’ (*i.e.*, twins that were brought up in different households per full-time custody bargains that divided the twins, but had claimed to have spoken regularly on the phone and often spent weekends together throughout their entire childhood). Given the above, a fair hypothesis would be that subjects living together develop common speaking habits (*e.g.*, phrases, frequently used words and jargon, and even accents). Thus, *nurture* plays its role, with hints that *nature* does as well. A recent study used computer vision and speech recognition technologies to verify kinship [137]. The authors showed significant gain in kinship verification performance when MM (*i.e.*, video-audio) is utilized using modern-day deep learning techniques that leverage both modalities, over just the still face images. In the end, Wu *et al.* used 400 video-pairs of *parent-child* pairs to show promise in the use of multi-modal systems for kinship recognition - which, to the best of our knowledge, has been the only attempt of using visual-audio data for recognizing kinship.<sup>1</sup>

Our contributions to the FR, biometric, anthropology, and MM communities were 3-fold.

- Built multimedia database: a large-scale dataset for kinship recognition was built using existing paired image data and an automatic labeling scheme. Media of different modalities is now available: video-tracks, audio segments, visual-audio clips, and text transcripts. Specifically, we extended the FIW imageset. Additionally, we restructured the MM family database to better encapsulate the added metadata and paired data respectively at the subject and instance levels.
- Recorded protocols and benchmarks: a new paradigm for kinship recognition suited for MM data as a step towards experimenting with real-world settings. Specifically, the problem has evolved from instance to template-based. Thus, we are the first to measure kinship recognition capabilities using a large-scale, multimodal template-based collection.

---

<sup>1</sup>The dataset collected for [137] is not available for public use.

## CHAPTER 6. FIW-MM



Figure 6.1: **Sample family of FIW-MM.** Top-to-bottom: *family-tree labels* show faces of members in the immediate family, with subjects of the same generations in the same row; *videos, audio, and contextual* exemplify sample video pairs of Dr. King Jr. and his daughter Andrea with tracklets of faces in the visual domain and audio data aligned frame-by-frame; *family photos* that contain Dr. Luther King Jr. randomly selected (note, cropped to fit); *faces* of Dr. King Jr. from adolescence-to-adulthood. Multiple faces are available for most subjects. Best viewed electronically.

- Showed the advantage of all modalities: Following the improved protocols and, thus, experimental settings, we demonstrate an increase in system performance from still-images, to still-images and videos, and then with audio speech signals added as well– a clear benefit of each added modality is shown. Our analysis highlights the shortcomings of the different media types for future work to address.

We believe this will attract a wider range of scholars to kin-based and multimedia problems. FIW-MM will be accessible online in various formats.<sup>2</sup>

## 6.2 Related Work

Early on, problems of recognizing kinship started with domesticated animals (*e.g.*, dogs [234] and sheep [235, 236]), as many species have a natural ability to recognize their kin through various signals (*e.g.*, touch, smell, visual, and acoustics). From this, we hypothesized that different types of media, besides image-level or conventional speech recognition, can be leveraged to better detect kinship in humans. Knowledge extracted from still-images and stationary speech signals are lacking an abundance of evidence. A more complex signal which helps improve decision making, such as dynamic features across video frames, can attribute inheritable characteristics (*e.g.*, expressions, mannerisms, and accents from different emotions). Nonetheless, such technology will take effort to acquire. We demonstrate the ability of the added modalities with face tracks from videos and standard audio features from speech signals.

We next review existing work in visual kinship recognition on still-images, and then more recent advances in the acquiring and modeling of visual-audio data for FR.

### 6.2.1 Kinship recognition

Computer vision researchers began using facial cues to recognize kinship about a decade ago. Specifically, Feng *et al.* proposed to model the geometry, color, and low-level visual descriptors extracted from faces to discriminate between KIN and NON-KIN [34]. Others then formulated the problem as various paradigms (*e.g.*, transfer subspace learning [106, 119], 3D face modeling [107], low-level feature descriptions [237], sparse encoding [9], metric learning [125], tri-subject verification [113], adversarial learning [184], ensemble learning [186], video understanding [211, 238, 239], and, most recently, video-audio understanding [137]). A common factor of the aforementioned

---

<sup>2</sup>FIW-MM - the data, code, trained models, and other resources - will be available upon publication of this work.

## CHAPTER 6. FIW-MM

was the attempt to improve discriminatory power for classifying a pair of faces as either KIN or NON-KIN; another commonality was the limited sample size and, thus, unrealistic experimental settings.

Robinson *et al.* introduced a large-scale image dataset to recognize families in still-imagery called FIW [4, 6]. FIW contains 1,000 families with an average of 13 family photos, 5 family members, and 26 faces. It came with benchmarks for 11 pairwise types, with the top performance of the baselines being a fine-tuned CNNs (*i.e.*, SphereFace [48] and Center-loss [158]). This was the beginning of big data in kin-based vision tasks— deep learning could then be used to overcome observed failure cases [11, 29]. Furthermore, new applications such as child appearance prediction [10, 240] and familial privacy protection [241] were done recently.

Nowadays, FIW continues to challenge researchers with various views of image-based tasks. A myriad of methods demonstrated the ability of machinery to use still-images to determine kinship in a pair or group of subjects. Nonetheless, only so much information can be extracted from still-images. The dynamics of faces in video data (*e.g.*, mannerisms expressed across frames) contain additional information, and audio as well as text transcripts (*i.e.*, contextual data describing the speech and other sounds) can widen the range of cues we model to discriminate between relatives and non-relatives. We propose the first large-scale multimedia dataset for kinship recognition. Specifically, we leveraged the familial data of the FIW image database to build upon the existing resource [4, 6], using the still-images of FIW and adding video, audio, audiovisual, and text data of subjects. Note that video, audio, and visual-audio differ in that the latter has the face speaking and the speech spoken are aligned, while the others are independent, unaligned clips. After its predecessor, we dubbed the database FIW-MM (Figure 6.1). En route to bridging research-and-reality, we follow the protocols of FIW [7], but now with the capacity to be template-based (*i.e.*, per National Institute of Standards and Technology (NIST) in [242]).

Besides the different use-cases, and independent research work it made possible, FIW was used as part of an annual data challenge motivated to attract more attention from and provide more incentive for the research community; namely, the RFIW series, which has been held each year since 2017 [5] to 2020 [25]. There have been many great attempts on the still-images as a result [243, 244]. Recent surveys [118], tutorials [30], and challenge papers [7, 112, 116, 120] elaborate on RFIW and the various submissions in detail.

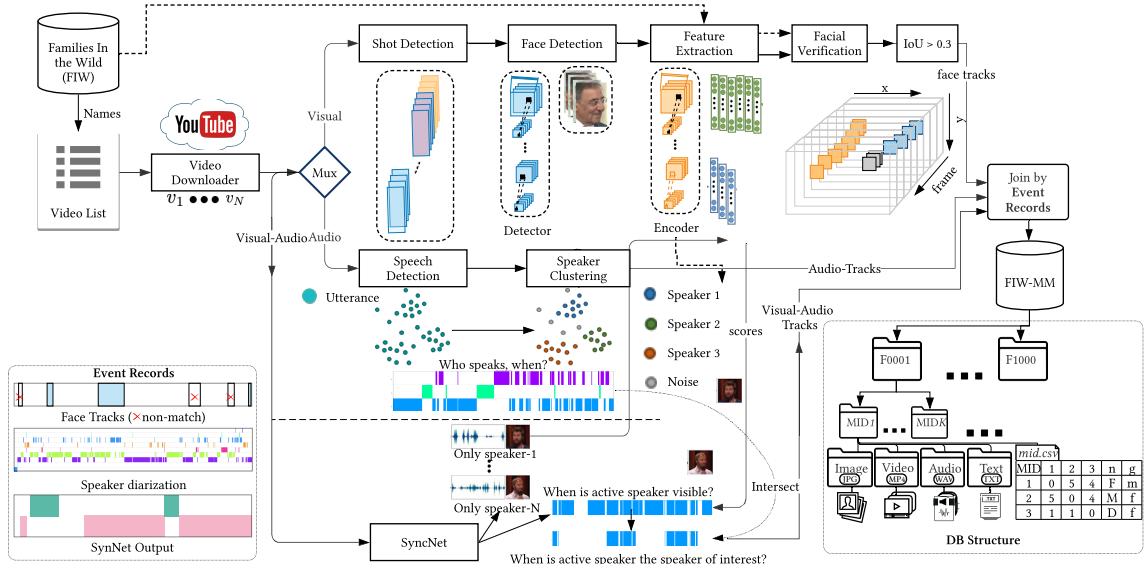


Figure 6.2: **Automated labeling framework.** For each of the 1,000 families, there are a set of  $K$  members. For this, the template of a member consists of all media available. Tag numbers 1-6 correspond to sections in Section 6.3.2

### 6.2.2 Audio-visual data

The archetypal big data resources for audio-visual identification problems are Voxceleb [224] and Voxceleb2 [225]. Similar to FIW-MM, the datasets were acquired by extending still-image collections (*i.e.*, Voxceleb and Voxceleb2 extended of the VGGFace [46] and VGGFace2 [14], respectively). Currently, the primary usage of Voxceleb is in speaker-based tasks, such as using the audio-visual data to detect and classify the speaker by the *who* and the *when* [223]. Additional speaker centric problems have been proposed using the Voxceleb collections, like to enhance speech signals [245], to detect *when* and *where* the speaking face is visible [246], and when the audio and mouth motions infer the lips and sound are in sync [247]. Nonetheless, the lip-reading task predates the larger Voxceleb with older lip-reading datasets [248, 249].

It is worth highlighting that these audio-visual databases were instrumental in applied research as well (*e.g.*, generating talking faces [250], where the input is a still-image face and a stream of audio, and the output are frames mocking the audio with the faces as if the input face was regurgitating the audio clip). In [251], face frames were generated from a still-image and audio clip, with pose information added as a control signal for the synthesized output. Furthermore, Voxceleb

predicted emotion labels via its own signals to automatically infer ground-truth [252].

Previous attempts to tackle kinship recognition have also been made with audio-visual data. Most relevant was in [137], where the authors built a collection made-up of 400 pairs. Wu *et al.* certainly demonstrated the core hypothesis of this work—multimedia can enhance our ability to automatically detect kinship in humans—as was clearly demonstrated in their work [137]. However, the sample size was limited in the number of pairs and in the types of labels, as there is no family tree structure, nor multiple samples per member (*i.e.*, age-varying), as is the case in our much larger and comprehensive FIW-MM.

### 6.3 The FIW-MM Database

The automated labeling pipeline in Figure 6.2 leveraged existing name labels available in FIW, and with each of these subjects represented by one-to-many face samples. Hence, the visual evidence of FIW was modeled and used to annotate the MM data by recording *when* and *where* events of interest occurred. In other words, we curated FIW by parsing videos to cropped face tracks, entire speaker instances, aligned visual-audio data, and the spoken words transcribed to text. All of these instances were organized as time-stamps representative of the start and end video frames, along with the bounding box information for the face tracks in all frames including and between the initial and final boundary frames. In this way overlap in samples was clearly identifiable. Furthermore, these instances were organized by family (*i.e.*, Family ID (FID)) and specific member (*i.e.*, Member ID (MID)) to allow possible outcomes to be limited to boolean (*i.e.*, *genuine* or *impostor* match). Family folders contained a folder for each of its MID, which held separate folders for face images, non-speaking face tracks, speech samples, transcribed conversations, and face tracks actively speaking (*i.e.*, visual-audio). Of course, relationship matrix and the genders for each MID remained in the family folder as was done back for the original FIW [4]. The structure of FIW-MM is shown bottom-right of Figure 6.2. Also, and along with being defined when introduced, various acronyms and symbols used in the following sections are listed under Acronyms at beginning of manuscript.

Next, the specifications, the pipeline used to automatically annotate the data, along with a few strategies used to further reduce the manual input requirements are described.

Table 6.1: **Database statistics.** Types are split based on the span in generation of the relationship.

	1 <sup>st</sup> -generation			2 <sup>nd</sup> -generation				3 <sup>rd</sup> -generation				4 <sup>th</sup> -generation				Total
	BB	SS	SIBS	FD	FS	MD	MS	GFGD	GFGS	GMGD	GGMGS	GGFGGD	GGFGGS	GGMGGD	GGMGGS	
# Subjects	883	824	1,542	1,914	1,954	1,892	2,041	426	463	483	526	39	30	45	37	13,099
# Families	345	334	472	666	676	665	670	154	174	178	191	9	10	11	10	953
# still-images	40,386	31,315	46,188	83,157	89,157	57,494	63,116	8,007	6,775	6,373	6,686	408	410	798	797	441,067
# Clips	123	79	81	155	134	147	138	16	18	25	15	2	4	0	0	937
# Pairs	641	621	1,138	1,151	1,253	1,177	1,207	263	280	292	324	28	18	36	28	8,457

### 6.3.1 Specifications

The goal was to extend FIW in the number of samples, the media types, and in the possible experimental settings. Recall that FIW provides name metadata and face images for an average of >13 members of 1,000 families [30]. Thus, we aimed to accumulate paired MM data for existing persons of FIW. Specifically, we acquired paired MM data for members of 150 families, with at least two members for each. Otherwise, if we only had samples for a single member for a family, the audio information could not be compared amongst any pair in the respective family. Hence, the requirement was set at least two members per family.

With complete access to FIW for research purposes<sup>3</sup>, we leveraged this data as the knowledge needed to build FIW-MM with minimal manual labor and zero financial cost. For this, we employed SOTA models and algorithms in speech and vision throughout the data pipeline. Our next steps are the modules and feedback loops making up the pipeline om Figure 6.2.

### 6.3.2 Data pipeline

Inspired by previous work, such as FIW [6] (*i.e.*, labeling families) and VoxCeleb [224] (*i.e.*, labeling audio-visual data), aspects of both were merged as the basis of our pipeline design and, in essence, one of three branches that make up our data collection pipeline. Specifically, the merging of the aforementioned pipelines make up the *audio-visual* branch, which processes end-to-end and in parallel with the *visual* and *audio* branches. The notion of branches is used for clarity in the following description, as each respective branch is concerned with the modality for which it is referred.

The following subsections cover the details of the pipeline built to acquire FIW-MM as the sequence of modules it grew to - the steps are covered in order of process (*i.e.*, from left-to-right in Figure 6.2). Philosophically, all data was assumed to be of type *non-match* (*i.e.*, zero amount of MM data to start). Then, there are various checkpoints throughout the branches that add data which

<sup>3</sup><https://web.northeastern.edu/smilelab/fiw/download.html>

was found to be a *match* with high confidence. Under the pretext that FIW-MM will be a resource used by experts from different data domains, all data points that *match* are saved (*i.e.*, visual tracks, audio, and audio-visual). Nonetheless, overlapping segments are clearly annotated such to remove repeated samples (*i.e.*, visual-audio will also be present in sets containing just visual and just audio). At the same time, if one or both modalities are of interest, then the maximum amount of data points is readily available. Note that no data points are repeated in sets created for included benchmarks. Also, the following subsections are numbered according to the yellow circle call-outs in Figure 6.2.

**1) Raw data resource.** FIW has still-image data for 1,000 families with over 13,000 family members (*i.e.*, subjects) in total. From the families, we chose a subset of 150 for which 2-5 members appeared in 1-3 YouTube videos, with a total of 500 subjects in 605 videos. The importance of this step was assuring that there were at least 2 members per family with MM; otherwise, the added modalities would have no basis to match about. Also, ethnicity for these 500 subjects were manually collected at this time. Video URLs were queried under unique Video IDs (VIDs) (*i.e.*,  $v_1, \dots, v_N$  for  $N$  videos). Generally speaking, the videos were either interview-style (*e.g.*, with a news anchor or alone in a plain room answering scripted questions) or face-time clips (*i.e.*, self-recordings of subject speaking directly to the camera, as is the normal case when face-timing).

Our scripts used Pypi’s youtube-dl<sup>4</sup> to download YouTube videos by URL, which were then archived under corresponding VID. Allowed with the MM (*i.e.*, in MKV file format), time-stamped captions were also scraped when available— later parsed as transcribed words spoken by the subject. Alongside the text, the MKV files were processed to three files: a copy of the original MKV for the *audio-visual* branch, and then an audio only (WAV) and visual only (MP4) extracted with *ffmpeg*. From the start, all video data was assumed a constant 25 FPS.

**2) Event records.** Before passing data down any branch, blank (sequential) tabular records were created for the duration of the video with tuples as index (*i.e.*, time and frame number)— one record per branch (*i.e.*, audio, visual, and audio-visual event records). These are essential for refinement processes that are later activated via a feedback mechanism. In essence, the mutual information across records at a given instance (*i.e.*, frame or time-stamp) are used to imply matches, contradiction, and non-matches across modalities (*i.e.*, a means to propagate labels across modalities). This usage of set theory helps both to validate true matches and filter out non-matches: others have too leveraged logic and sets to parse videos [253]; however, opposed to high-level semantics such as types of objects present, we reference output of simpler tasks (*e.g.*, face or no face, speech or non-speech, same or

---

<sup>4</sup><https://github.com/yt-dlp/yt-dlp/>

Table 6.2: **Task-specific counts.** For individuals (**I**), families (**F**), still-face images (**S**), video-clips (**V**), audio snippets (**A**), audio snippets (**VA**) in the set of probes (**P**), gallery (**G**), and in total (**T**).

	Train					Val					Test				
	I	F	S	V	A	I	F	S	V	A	I	F	S	V	A
$T^I$	T 2,976	571	16,464	290	7,217	955	190	5,458	72	3,308	972	192	5,231	91	1,775
$P$	571	571	3,039	47	1,843	190	190	1,334	16	789	192	192	993	23	876
$G$	2,475	571	13,571	244	5,581	791	190	4,538	56	2,519	800	192	4,705	69	899
$T$	3,046	571	16,610	291	7,424	981	190	5,872	72	3,308	992	192	5,698	92	1,775

different face or voice)– this increases random chance and thus reduces low confident decisions.

**3) Visual branch.** We first split a video into scenes using two global measures under the assumption that, statistically, neighboring frames will match as close as 90%: HSV (*i.e.*, color) and local binary patterns [156] (*i.e.*, texture) features were extracted and used to parameterize two probabilistic representations per frame, which were compared using KL-Divergence and compared to a threshold of 0.1 [254]. This produces a set of shots for each of the  $V$  videos of size  $C$ , *i.e.*,  $v_c \in \{1, \dots, C\}$  represents all shots detected in the  $i$ -th video. From there the first, last and the frame in between closest to the centroid (in color and texture) of the entire track (*i.e.*, the beginning, end, and the assumed best representation for the respective clip). The three frames per clip are then passed through a MTCNN face detector [102], and clips with no faces detected in at least one of these frames. In addition, the set of clips is filtered further by comparing detected faces to the ground-truth faces of FIW. Again, clips with no matches are discarded. Note that this was a means to quickly drop unwanted data. To compare faces, faces were encoded with ArcFace via the architecture, settings, training details, and *matcher* in [49]. Specifically,

$$d_{bool}(x_i, x_j) = d(x_i, x_j) \leq \theta, \quad (6.1)$$

where the *matcher*  $d_{boolean}$  compared the  $i$ -th face detected to  $j$ -th FIW face encoding  $x$  [150]. In other words,  $d_{boolean}$  is the decision boundary in similarity score (or metric distance) space— if threshold  $\theta$  is satisfied, assume match; else, non-match. Note that it is currently assumed that  $i$  and  $j$  are from different sets (*i.e.*, with  $J$  labeled samples of a subject from FIW and  $I$  face detections in the new video data). The *matcher* in Eq 6.1 was set as cosine similarity the closeness of the L2

## CHAPTER 6. FIW-MM

normalized [255] encodings by  $d_{bool}(x_i, x_j) = 1 - d(x_i, x_j) = \frac{z_i \cdot z_j}{\|z_i\|_2 \|z_j\|_2} > \theta$ , where  $z$  represents an encoded piece of media. At this stage,  $\theta = 0.2$  was manually set for a high recall. The matching process - including the usage of ArcFace to encode faces - is the *matcher* used throughout.

Next, the MTCNN outputs were generated for all frames in clips, while saving the bounding box coordinates, fiducials (*i.e.*, 5 points), and confidence scores. Next, only continuous face tracks in clips were kept. For this, the ROI was set on the previous location of the face, and then IoU was calculated frame-by-frame, each value must surpass a threshold of 0.3. Finally, up to 25 faces were sampled uniformly from track (*i.e.*, opposed to choosing the top  $K$  based on pose information, as this yielded redundancy in similar frontal posed faces). Each was then passed to  $d_{bool}$  with each of the  $I$  labeled faces (*i.e.*, producing  $K \times I$  score matrix). The mean across  $I$  samples was calculated to produce a single score per the  $K$  faces, at which point the value at the 25-percentile was compared to  $\theta = 0.25$ . The fusion of scores was done in such a way to both consider all the existing labeled faces equally, while avoiding a few (of the  $K$ ) low-quality detections having any weight. Upon this process, and with the aid of SOTA techniques mentioned throughout, this step alone yielded many face tracks matching with a high confidence.

**4) Audio branch.** Audio data, in its raw form, is extracted from the videos and saved as high-quality wave files. We first set out to do speaker diarization on each video: we aimed to have a record indicating the presence of speech, from which change in speaker is marked, and, ultimately, the number of speakers in the video along with *who* speaks *when*. Note, we assume no audio labels. Thus, the speakers are arbitrarily tagged per video.

Put differently, the first purpose of this branch is to find the number of speakers per video, with predictions based on the detected speakers on who spoke when: a speech detector determined the *when*, and then clusters all the different speech segments to determine the number of speakers and, thus, which speech segment to assign to which of the speakers (*i.e.*, the *who*). The former was implemented using PyPi’s SpeechRecognizer<sup>5</sup>, with the latter based on models from [256]. See supplemental for further detail. Finally, parsing through segments and marking as speaker<sub>a</sub>, speaker<sub>b</sub>, ..., speaker<sub>j</sub>, where  $j$  is the number of speakers in a given clip. The time-stamps are used to detect speakers of interest.

**5) Visual-audio branch.** Focused is on detecting when the speaker is in the field of view. Thus, its purpose was to detect the boundaries (*i.e.*, start and end frames) for which the face and speech are aligned. An intuitive way to do this is to relate the faces detected and the lip movement with the

---

<sup>5</sup>[https://github.com/Uberi/speech\\_recognition](https://github.com/Uberi/speech_recognition)

Table 6.3: **TAR at specific FAR.** Scores are for template-based settings: still-images only (left column), +videos (middle), and +video+audio (right). Higher is better.

FAR/TAR (%)	BB			SS			SIBS			FD			FS			MD			MS			Average		
0.5 (EER)	97.8	97.8	<b>97.8</b>	91.5	92.3	<b>92.3</b>	91.7	90.8	<b>90.8</b>	79.8	77.8	77.8	85.3	85.3	<b>85.3</b>	90.6	88.8	88.8	81.3	82.6	<b>82.6</b>	88.3	87.9	<b>87.9</b>
0.3	94.1	94.1	94.1	88.0	87.2	87.2	82.9	83.9	83.9	63.5	66.5	66.5	77.1	79.1	79.1	82.4	82.0	82.0	68.9	70.1	70.1	79.6	80.4	<b>80.4</b>
0.1	88.1	87.4	87.4	76.1	76.1	76.1	68.7	68.2	68.2	34.5	36.9	36.9	54.3	54.3	54.3	62.2	63.1	63.1	46.1	46.5	46.5	61.4	61.8	<b>61.8</b>
0.01	70.4	70.4	70.4	54.7	55.6	55.6	44.2	46.1	46.1	5.9	7.9	7.9	23.6	24.0	24.0	28.3	31.3	31.3	11.6	13.3	13.3	34.1	35.5	<b>35.5</b>
0.001	54.8	<b>57.0</b>	<b>57.0</b>	47.9	<b>48.7</b>	<b>48.7</b>	<b>29.5</b>	29.0	29.0	2.0	<b>2.5</b>	<b>2.5</b>	9.3	<b>10.9</b>	<b>10.9</b>	14.2	<b>14.6</b>	<b>14.6</b>	3.3	<b>4.6</b>	<b>4.6</b>	23.0	<b>23.9</b>	<b>23.9</b>

audio— which is at the core of many speaker identification methods in MM [257]. To acquire this, videos were processed using SyncNet [258]) with the settings and trained weights from [259]. The output were tracks: first trimming the video about the boundaries detected, and then cropping out the faces using the detected bounding box coordinates extended 130% in all four directions. From this, each track is static spatially, and with each face detection captured within. This modification made it so individual tracks were of constant size and location in pixel space; opposed to producing tracks with moving coordinates to preserve the face in the field of view (*i.e.*, the added 30% covered this). We then had three sets of coordinates saved (*i.e.*, the original detection, the extended version, and the set accounting for relative offsets for the crop). Similar to the *visual branch*, labeled faces from FIW were then used to determine the tracks belonging to the subject of interest. Once compared and, thus, filtered, all cropped tracks were manually inspected. Upon this the paired data was assumed.

## 6.4 Problem Definitions and Protocols

The FIW-MM database extends the large-scale imageset FIW [4, 6]. Specifically, the images and names of FIW, as explained in the previous section, allowed labeled multimedia data to be acquired via an automated process. Following the protocols of the recent RFIW data challenge [7], we benchmark two kin-based tasks: verification and search & retrieval. Differences from FIW and, thus, from the experimental settings of RFIW, are protocols based on still-imagery— uni-modal and the experiments are organized as *one-shot* problems. In contrast, FIW-MM offers multiple modalities, resulting in many more samples and sample types (Table 6.1). Furthermore, in an attempt to further bridge the gap between research-and-reality, the protocols we explain next is the first attempt in kinship recognition to follow template-based protocols [242].

As for experimental tasks, kinship verification has been the primary focus. More recently, the emergence of the more challenging but more practically awarding task of *searching for missing*

## CHAPTER 6. FIW-MM

*children* task [7], *i.e.*, search and retrieval. We benchmark FIW-MM for both these tasks. However, opposed to the single-shot setting followed up until now, we use templates [242]; hence, a means to move experiments closer to settings for operational use-cases.

Template-based experiments are organized as follows: Known subjects (*i.e.*, prior knowledge of identity and family) are first enrolled in a *gallery*  $G$ . At inference, the aim in the search and retrieval task is to compare an unseen *probe*  $P$  to subjects enrolled in  $G$ ; the verification task compares a list of *probes* to individual *gallery* subjects (*i.e.*, one-to-one) with the solution space of either *KIN* or *NON-KIN*; kinship identification compares the *probe* to the entire *gallery* (*i.e.*, one-to-many), with the end result being a ranked list of family members. In all cases, at least one family member exists in  $G$ , making for a closed-set recognition problem.

Specifically, template  $X$  holds all of the media for a subject (*i.e.*, face images, videos, audio-clips, and text transcripts). Hence,  $X$  consists of samples  $x$ , where each  $x$  is an independent piece of media represented as an encoding  $z$ . For instance, still-image  $x$  encoded as  $z$  by  $\mathcal{F}(x) = z$ , where  $\mathcal{F}$  maps faces to a learned feature space (*i.e.*,  $\mathcal{F}(x) \in \mathbb{R}^d$ , where the dimension  $d$  represents the size of the respective encoding). The same is done for face tracks in videos, which were fused to a single encoding by average pooling. Put formally, a face track is represented as  $\bar{z} = \frac{1}{m} \sum_x \mathcal{F}(x)$ , where  $m$  is the frame count. Similarly, an audio segment (*i.e.*, a clip where subject speaks without interruptions or major pauses) is treated as a single piece of media  $x$  via average pooling all encodings to form a single representation per clip. Note that a video may consist of several independent visual, audio, and visual-audio (*i.e.*, aligned) tracks. Thus, there are many independent media samples for both the visual and audio modalities. Again, subjects are represented by these templates  $X$  made up of these various media samples  $x$ , such that the  $j^{th}$  subject can be represented by  $k$  media samples as follows:  $X_j = \mathcal{F}_t(x_1), \mathcal{F}_t(x_2), \dots, \mathcal{F}_t(x_k)$ , where  $t$  corresponds to the media type and, hence, the corresponding encoder. From this,  $|X_j|$  is the total number of encodings for subject  $j$ . The *gallery*  $G$  consists of a set of subjects by  $G = \{(X_1, y_1)^l, (X_2, y_2)^l, \dots, (X_n, y_n)^l\}$ , where  $y$  are identity labels for each of the  $N$  subjects, and  $l \in \{1, 2, \dots, L\}$  are ground-truth for  $L$  families. To establish a precise definition for problems of kinship, each tuple also contains a tag representing the set of  $L$  families (*i.e.*,  $(X_j, y_j)^l$ ), where  $l \in \{1, 2, \dots, L\}$ . Further partitioning of the data is done per requirements of a task. For instance, for the verification, the  $m^{th}$  pair of tuples from the same family  $\mathbb{P}_m = ((X_i, y_i) \cap (X_j, y_j))$ , where  $i \neq j$ , inherit labels *KIN* (*i.e.*, *match*) and relationship type.

Following the 2020 RFIW, each task consists of a train, validation, and test set. These sets are disjoint in family and subject IDs, and are roughly split 60%, 20%, and 20% for the train, validation, and test set, respectively. Thus, the splitting is done using the family labels, and the

resulting partitioning of sets is static for all tasks.

### 6.4.1 Kinship verification

Kinship verification is a challenging task within a complex topic. It inherits all the challenges of traditional FR, with aspects amplified in difficulty due to kinship being a soft attribute with high variation, bias in nature, and directional in the variety of relationship types. The most fundamental question asked in kinship verification, and re-asked in all other kinship discrimination tasks is whether a face pair is related. Therefore, kinship verification is a boolean classification of pairs (*i.e.*,  $y \in \{KIN \cup NON-KIN\}$ ). Knowledge of the relationship type is assumed to be known. Thus, provided the output of the model for a given pair is *KIN*, then the specific type is implied. Future efforts could incorporate relationship-type signals to advance capabilities of kinship detection systems; however, and as stated upfront, verification provides the simplest of all the benchmarks and, up until now, is the most popular [7].

#### 6.4.1.1 Data splits and settings

The data is organized as pairs, with pairs a part of a set of common relationship-type. Specifically, pairs are of type BB, SS, or SIBS of mixed-sex (*i.e.*, same generation), or FD, FS, MD, or MS (*i.e.*, difference on 1-generation). Counts for all types of relationship pairs are listed in Table 6.1, with the aforementioned types (*i.e.*, same and 1-generation) used in experiments provided sample sizes are such to allow for fair representations. Data splits (*i.e.*, train, validation, and test) and their sample counts are listed in Table 6.2. The task here has no concept of query and gallery.

#### 6.4.1.2 Metrics

The one-to-one paradigm (*i.e.*, kinship verification) is the main view vision researchers aim to solve. The task is to determine whether a face-pair are blood relatives (*i.e.*, *true kin*). Conventionally, a query consists of a single face image  $x_1$ , which is then paired with a second face  $x_2$  to predict against (*i.e.*, a one-shot, boolean classification problem with labels  $y \in \{KIN, NON-KIN\}$ ). Put formally, given a set of face-pairs  $(x_1, x_2)_s^m$ , where the number of sample pairs  $s \in \{1, 2, \dots, S\}$  of relationship-type  $m \in M \rightarrow \{BB, SS, \dots, GMGD, GMGS\}$  (*i.e.*,  $|M| = 11$ ). A set of pair-lists  $\mathbb{P} = \{[(x_1, x_2)_1^m], [(x_1, x_2)_2^m], \dots, [(x_1, x_2)_S^m]\}_1$  for the  $M$  types, and with the label determined by the indicator function  $\mathbb{1}$  for a single pair  $\mathbb{P}_s \rightarrow \{0, 1\}$ , *i.e.*,

Table 6.4: **Identification results, with TAs highlighted.** Accuracy scores for different ranks are listed (*i.e.*, higher is better). Also, MAP scores are provided for each.

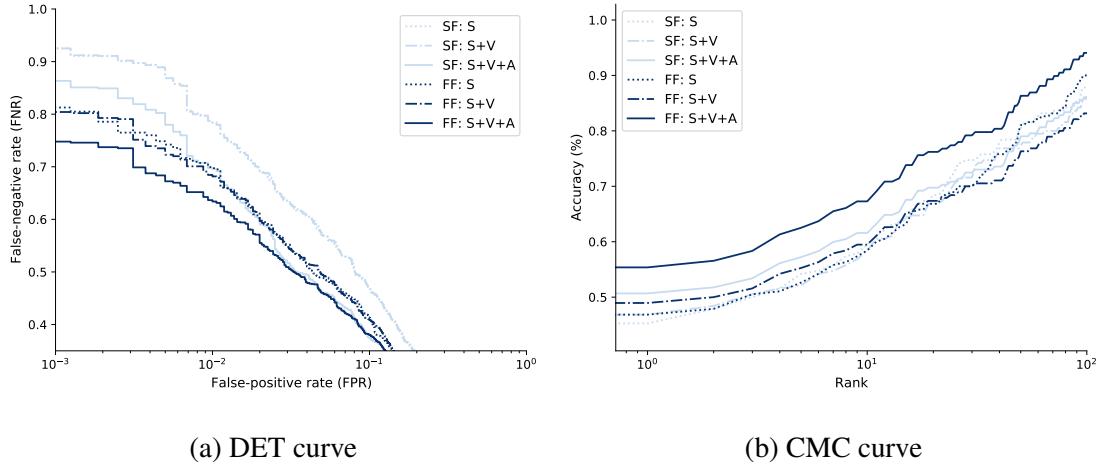
		Rank					
		@1	@5	@10	@20	@50	mAP
<i>img</i>	mean	0.29	0.43	0.54	0.64	0.78	0.13
	median	0.28	0.44	0.52	0.64	0.77	0.13
	max	0.11	0.19	0.28	0.34	0.52	0.06
	TA	0.31	0.43	0.52	0.63	0.74	0.14
<i>img+video</i>	mean	0.30	0.44	0.52	0.64	0.77	0.14
	median	0.28	0.44	0.50	0.63	0.76	0.14
	max	0.13	0.21	0.26	0.30	0.44	0.06
	TA	0.34	0.46	0.55	0.68	0.75	0.16
<i>img+video+audio</i>	mean	0.30	0.44	0.52	0.64	0.77	0.14
	median	0.28	0.44	0.50	0.63	0.76	0.14
	max	0.13	0.21	0.26	0.30	0.44	0.06
	TA	<b>0.56</b>	<b>0.59</b>	<b>0.63</b>	<b>0.74</b>	<b>0.78</b>	<b>0.24</b>

$$\mathbb{1}(\mathbb{P}_s) = \begin{cases} 0 & \text{NON-KIN} \\ 1 & \text{KIN} \end{cases}. \quad (6.2)$$

Note, a  $\mathbb{P}_s$  consists of a pair of templates and, thus, the task is to determine whether the media of the templates provide evidence of the two subjects being blood relatives; notice Eq. (6.2) is the template *matcher* defined in Eq. (6.1).

As described, FIW-MM is organized as templates with many samples from various modalities (*i.e.*, still-face, face-tracks, audio, and transcripts (contextual). Specifically, true IDs  $y$  are paired with a template of all media available for the respective subject. In contrast with conventional kinship recognition, where one image is compared to another, the one-to-one paradigm is based on templates (*i.e.*, one template is compared to another). For consistency, given  $\mathbb{P}_s^m = ((X_i, y_i), (X_j, y_j))$  as a pair of templates for different subjects (*i.e.*,  $X_i$  and  $X_j$ , where  $i \neq j$ ).

DET curves, along with average verification accuracy, were used for kinship verification as also were TAR across intervals of FAR (Table 6.3).



**Figure 6.3: Plotted results.** Shown are score fusion (SF) and feature fusion (FF) as late and early fusion methods, respectively. Included are still-images  $S$ , video clips  $V$ , and audio segments  $A$ , with still-images and video were fused  $S + V$ , and also still-images, audio, and video were fused  $S + V + A$ . Clearly, both tasks benefit from early fusion: the DET curve (left) summarizes the verification task by plotting FNR as a function of FPR (*i.e.*, lower is better); search and retrieval is summarized as a CMC (right) by showing the accuracy as a function of rank (*i.e.*, higher is better).

## 6.4.2 Search & retrieval (missing child)

### 6.4.2.1 Overview

Kinship identification is organized as a *many-to-many* search and retrieval task, with each subject having one-to-many media samples. Thus, we imitate template-based evaluation protocols [242]. Furthermore, the goal is to find relatives of search subjects (*i.e.*, *probes*) in a search pool (*i.e.*, *gallery*).

### 6.4.2.2 Data splits and settings

A gallery  $G = \{g_i\}, (i = 1, \dots, N)$  is queried by a set of probes  $P = \{p_j\}, (j = 1, \dots, M)$  for search and retrieval, where  $g_i$  is the  $i$ -th template in  $G$  and  $p_j$  is the template of the  $j$ -th query subject. As mentioned, a template consists of samples of various modalities. Given a template of MM, various schemes were applied to integrate the ID information from all media components of  $P$ .

### 6.4.2.3 Metrics

Scores of  $N$  missing children are calculated as

$$AP(l) = \frac{1}{P_L} \sum_{tp=1}^{P_L} Prec(tp) = \frac{1}{P_L} \sum_{tp=1}^{P_L} \frac{tp}{rank(tp)},$$

where average precision (AP) is a function of family  $l \in L$  (*i.e.*,  $|L| = P_L$ ) for a given TPR. All AP scores are averaged to find the mean AP (*i.e.*, MAP):

$$MAP = \frac{1}{N} \sum_{l=1}^N AP(l).$$

Also, Cumulative Matching Characteristic (CMC) curves as a function of rank enable for analysis between different attempts [260], along with the accuracy at rank 1, 5, and 10.

## 6.5 Benchmarks

### 6.5.1 Methodology

The problems of FIW-MM have various views— multi-source and multi-modal. The former varies in samples and in treating the different media-types independently until the matching function outputs scores (*i.e.*, late-fusion). The latter demands a method for early fusion (*e.g.*, feature-level) which should enhance performance by leveraging informative samples while ignoring noisy and less discriminative samples. We next describe the modality-specific features (*i.e.*, encoding different media types), and early fusion.

#### 6.5.1.1 Visual features

FR performance traditionally focuses on verification— popularized by the Labeled Faces in the Wild dataset [150] (images) and the YouTubeFaces dataset [165] (videos). In contrast, the newer IJB-[A,B,C] FR datasets [242] unifies evaluation of one-to-many face identification with one-to-one face verification over templates (*i.e.*, sets of imagery and videos for a subject). Then, visual kinship recognition research followed a similar path, addressing the simpler verification task. FIW-MM provides the data needed to run template-based kin recognition experiments.

We demonstrate results from a variety of naive fusion techniques (*e.g.*, average pooling of features or voting of scores). To no surprise, the score-based fusion outperforms the naive feature-

level fusion schemes. Specifically, the mean of all scores, both within a template and comparing templates (Table 6.3). The gain from each added modality is clear from just the naive score-fusion.

As mentioned, naive fusion methods at the feature level are an ineffective way of combining knowledge. Provided a collection of media - media that varies in modality, quality and discriminative power - a simple, unweighted average across the items of a template does not exploit all available information. To better fuse the template, we adapt a model to the template to best represent the subject for verification or identification of family members. Details are provided right after the description of audio features.

### 6.5.1.2 Audio features

All speech segments were encoded a SOTA deep learning architecture [256]. Specifically, we trained SqueezeNet [218] as a 34-layer ResNet [16] with an *angular prototypical loss* and optimized with Adam [261] to transform WAV-encoded audio files to a single encoding, *i.e.*,  $f(x) = z \in \mathcal{R}^d$ , where  $d = 512$ . *Angular prototypical loss* [262] learns a metric alongside softmax to minimize within-class scatter (*i.e.*, penalty formed as the sum of euclidean distances from all samples of a subject in a mini-batch from the mean centroid of the respective mini-batch). Specifically, a support set  $S$  and a query  $Q$  are set in each mini-batch on a subject-by-subject basis, with  $Q$  made-up of a single utterance to compare with the centroid of  $S$  that consists of all other samples in the mini-batch for that class. *Angular prototypical* takes advantage of the perks of using centroid prototypes, while enhancing by following generalised end-to-end (GE2E) [263] usage of a cosine-based similarity metric. This is scale invariant, is more robust to feature variance, and facilitates stability in convergence during training [264].

### 6.5.1.3 Feature fusion

TA [178] is a form of transfer learning that fuses the deep encodings of many labeled faces from a source domain with a template specific SVMs trained on the target domain. For kinship verification, we employ *probe adaptation*; for identification (*i.e.*, search & retrieval), we do *gallery adaptation*. Thus, TA enabled early fusion of different media types in both tasks.

Specifically, a similarity function  $s(P, Q)$ , for probe  $P$  and reference template  $Q$ , is learned for a given probe (*i.e.*, template). An SVM is trained on top of the face encodings with media in  $x^+$  as the positive samples and the set of negatives  $x^-$  being single sample from subjects in the train set (*i.e.*,  $x^+ \ll x^-$ ). For verification, this process repeats for another SVM  $Q$  (*i.e.*, the template of

## CHAPTER 6. FIW-MM

the subject in question). Negatives were set in same way. Then, let  $P(q)$  represent the evaluation of media encodings of template  $Q$  upon being trained on  $P$ . We do this in both direction via

$$s(P, Q) = \frac{1}{2}P(q) + \frac{1}{2}Q(p). \quad (6.3)$$

The score produced is the result of the templates fused together from media to an SVM and then to a score.

The benefit of SVMs is in the kernel. Specifically, the linear, max-margin modeling scheme of a vanilla SVM has proven effective at separating non-linear feature space of two classes; (*i.e.*,  $i$  and  $j$ , where  $y_{ij} = \pm 1$  for instances of the same (+) and different (−) classes. Thus, the implicit embedding function (*i.e.*, kernel)  $K(x_i, x_j, y_{ij}) = \varphi(x_i, y_i)\varphi(x_j, y_j)$  projects the encoding pair to a non-linear space such that the SVM learns the best hyperplane  $\mathbf{w}^T K(x_i, x_j, y_{ij}) + b = 0$  to separate the two classes. This is done by (1) maximizing the margin and (2) minimizing the loss on the training set—weights  $\mathbf{w}$  is learned, while bias term  $b$  we set to 1 (*i.e.*, concatenated on  $\mathbf{w}$  as an added dimension). Also,  $K(x_i, x_j, y_{ij}) = \exp \frac{\|x_i - x_j\|^2}{2\sigma^2}$  for  $y_{ij} \in \{-, +\}$  as the respective class (*i.e.*, Gaussian RBF kernel [265] projects all encodings to a higher dimensional space); then the predicted class is inferred as  $\hat{y} = \mathbf{w}^T \varphi(x_i)\varphi(x_j) + b$ . We used *dlib*'s [266]—L2 regularized cosine-loss with class-weighted hinge-loss, *i.e.*,

$$\begin{aligned} \min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \lambda_+ \sum_{i=1}^{N_+} \max [0, 1 - y_i \mathbf{w}^T f(x_i)]^2 \\ + \lambda_- \sum_{j=1}^{N_-} \max [0, 1 - y_j \mathbf{w}^T f(x_j)]^2. \end{aligned} \quad (6.4)$$

Adapting this to the notion of a gallery, the protocols are set for *gallery adaptation*: train a similarity function  $s(P, G)$  from a probe  $P$  to gallery  $G$ . A gallery of templates  $G = \{X_1, X_2, \dots, X_m\}$  are used to train the SVM (*i.e.*, the scoring function  $s(P, X_i)$ ). The difference between *probe adaptation* and *gallery adaptation* is in the negative sets. Along with the sample per subject trained against for *probe adaptation*, *global adaptation* samples all other templates in  $G$  as additional negatives. Again,  $N_+ \ll N_-$ . The class imbalance is handled via class-weighted hinge-loss in Eq 6.4, with  $\lambda_+ = \lambda \frac{N_+ + N_-}{2N_+}$ ,  $\lambda_- = \lambda \frac{N_+ + N_-}{2N_-}$ , which are regularization constants inversely proportionate to class frequency. The constant  $\lambda$  trades-off between the regularization and loss, which we set to 10 as in previous work [178].



Figure 6.4: **Random hard sample.** Template of a true FS pair that was incorrectly classified using score fusion, but correct for TA (*i.e.*, feature fusion). Here only a single face is available for the father (left), while all instances of the son are at a young age (right).

### 6.5.2 Results

The ability of a system to discriminate is improved with each added modality (Fig 6.3, Table 6.3 and 6.4). Considering the benchmarks use conventional speech and FR technology, and our hypothesis that video and audio boosts discrimination, much promise reflects—these notable improvements would likely continue to climb provided a more sophisticated or specific solution. It would be interesting to fuse earlier on than done here, and train machinery jointly with audio-visual data. This way, more complex dynamics of facial appearance, along with the corresponding sound of voice, could further improve and give additional insights.

There is a trend in the type of samples that were corrected when comparing the score fused to the feature fused (*i.e.*, TA) results. As shown in Figure 6.4, the challenge of recognizing kinship from samples of one or more member at a young age is mitigated. TA learns to better discriminate in such conventional failure cases. Additionally, some templates made up of multiple instances, often are better than others when comparing. Hence, TA does not simply average all instances with equal weighting, as done in late (*i.e.*, score) fusion—seen in cases containing a minority of samples that are more discriminative than its majority (Figure 6.5).

### 6.5.3 Discussion

The template-based protocol adds practical value by mimicking the more likely structure posed in operational settings, per NIST [242]. Besides, several other factors make it a more interesting formulation and therefore, a higher potential for researchers to show-off their creativity. For instance, opposed to using a single sample per subject (*i.e.*, one-shot learning), each now is represented in a

set of media (*i.e.*, a template). The questions now arise - how to best fuse knowledge and incorporate evidence from different modalities, and how to best learn from all available MM data? Another consequence of using templates is that the random chance is increased from (1) the knowledge added to pool (or fuse) from the added modalities, and (2) the gallery size reduces from tens of thousands by nearly ten-fold. The latter is not an implication of lesser difficulty, but the byproduct of reducing bias in data [25]. That is, opposed to having one-to-many samples per subject, there is just one template. Mitigating certain sources of data imbalance (*i.e.*, whether there are thirty samples or just one) a system's ability to recognize a particular pairing or group affects the metric evenly for all. In other words, a system may easily recognize a specific parent-child pair - regardless of the number of face samples and, consequently, the number of face pairs. Hence, the impact on the metric is proportional to the number of unique pairs, not sample pairs.

## **6.6 Future Work**

FIW-MM pushes the bar for possibilities in automatic kinship recognition and understanding. One immediate next step for research involves the benefit of gathering experts of different domains, such as those in sequence-to-sequence modeling, whether visual (*i.e.*, video), audio (*i.e.*, speech), contextual (*e.g.*, conversations, parts-of-speech, etc.), or early-fusing pairs or groups. Let us next discuss a variety of ways the data is foreseen to benefit and bring together different research communities, and beyond (*i.e.*, its inherent commercial potential).

We expect FIW-MM will bring experts of anthropology and genealogy-based together with those researching MM, machine learning, and vision topics toward helping to identify the hidden patterns that relate families in the MM data. Particularly, let us consider audio. As we have shown, pre-trained models from the speech recognition domain provide a means to acquire audio features with discriminative power that boosts kinship recognition systems which use only visual evidence. Furthermore, high-level semantics (*i.e.*, attributes) like accents, commonly used phrases, and speaker demeanor, could boost the performance of a system and also provide insights by interpretation. Similarly, studies focused on familial language components and inherited changes; or even the same generation (*i.e.*, commonalities and differences in the spoken tendencies of siblings) can be quite revealing as well. Hence, the potential only grows with audio-visual data, both in model complexity (*i.e.*, capacity) and practical uses: one could model mannerisms from the dynamics captured across video frames, and provide the answers to questions like “do I have my mother’s smile?”.

The data mining potential of FIW-MM is noteworthy. The family trees, abundance of data



Figure 6.5: **Hard sample pair.** A true MS pair incorrectly classified using late fusion, while correctly identified as KIN when using the early fusion via TA. Challenges here are in the young age of the son and the majority of the faces of the mother occluded by sunglasses (*i.e.*, score fusion puts equal weight on all samples, where TA learns to better discriminate).

points, rich metadata for individuals and relationships among MM data—FIW-MM could serve as a basis for group-based (*i.e.*, social) data mining. Additional data can enhance or target specific nature-based studies, traditional ML-based audio, visual, and audio-visual tasks, or even further extend this dataset. Fusing audio-visual data is an ongoing, unanswered problem [231]. Note the following: (1) the model training, for instance, with one or multiple incomplete modalities, (2) the data processing and balancing, and (3) the underlying roots of the problem to the high-level semantics, similar to contemporary biometrics systems with audio-visual data—FIW-MM and, thus, this work in its entirety, poses more problems than it solves. We actually introduce a much larger problem space than that of solutions.

Note that FIW-MM and thus this work in its entirety pose more problems than it solves: from the model training, to improvements made when dealing with incomplete modalities, and even the data processing and data imbalance; from the underlying roots of the problem to the high-level semantics, similar to contemporary biometrics systems with audio-visual data; we introduce a much larger problem space than that of solutions.

Another direction is fusion. For experiments, we included early and late fusion by joining the different media as features and scores, respectively. Scores were fused naively, ignoring the signal type, and assuming all samples and media types should be weighted uniformly. Fusion can incorporate more sophisticated techniques: cross-modality, selectively choosing the highest quality samples, or a decision tree based on modalities. This concept alone is vast in empty solution space—

## CHAPTER 6. FIW-MM

whether data fusion, where the input is then clips of aligned audio-visual data; early-fusion, which was exemplified with TA fusing the features; or late-fusion, also demonstrated by averaging scores, but could have just as easily been guided by a more clever decision tree mechanism. Besides, meta-knowledge, like relationship types (*e.g.*, directional relationships that inherently exist), genders, age, and other attributes, could indicate final decisions. Hence, there are an abundance of fusion paradigms—none are trivial, yet most hold promise.

Research topics to spawn off the proposed is vast, to say the least; the specifics suggested here are limited by our perception. We expect scholars and experts of different domains to seek out paradigms not thought of by us in the moment. Hence, whether it be an improved variant of adapting templates and feature fusion (*e.g.*, like in [267]), deciding when to fuse, a new method of integration, along with the integration details, are all open research questions.

In the end, the data resource outweighs the benchmarks. This is by design, as this resource will be readily available for research purposes - even a complete characterization of the contents as is (*i.e.*, ablation studies like on the effects of template sizes, media type versus relationship types, or even high-level interpretations (*e.g.*, smiling faces versus neutral)).

## 6.7 Conclusion

We introduced new paradigms (*i.e.*, template-based) for kinship recognition via the proposed FIW-MM database. FIW-MM contains audio, video, audio-visual, and text captions for 2+ members from 150 / 1,000 families of FIW. Our labeling pipeline uses multi-modal evidence and a simple feedback schema to leverage the labeled data of FIW to propagate ground truth for the added modalities. Benchmarks show improved performance with each added media type, and then further by early fusion. FIW-MM marks a major milestone for kin-based problems by welcoming experts of other data domains. In addition, FIW-MM supports a number of MM recognition tasks due to its rich metadata, template-based structure and multiple modalities.

One motivation of this survey is to establish cohesive views of the major milestones via protocols that are clearly defined, data splits that are ready for download, and trained models that make baselines reproducible.<sup>6</sup> Hence, components to reproduce experiments that we report make up part of the supplemental material. We cover the edge cases that challenge SOTA, including an examination of the different settings and training tricks that further our abilities in kin-based detection from faces.

---

<sup>6</sup><https://github.com/visionjo/pykinship>

## **Part III**

# **Post Processing**

# **Chapter 7**

## **Kinship Recognition - State of Technology**

### **7.1 Overview**

To review the current state of technology accessible for automatic kinship recognition in multimedia there are two separate aspects of the problem in need of elaboration: our limitations and the challenges for which they are set and details of the real-world uses-cases that have been mentioned briefly throughout this dissertation.

Hence, having reviewed the means and the results, let us now examine how it fits in practice. Specifically, let us now summarize the technical challenges still relevant in problems of automatic kinship recognition. We cover the challenges as they exist: general challenges as seen collectively, the way in which the challenges put limits on capabilities for current SOTA machinery, and in which ways challenges of the problem are set by nature, the environment, and inherited by the source and structure of the data in itself. We then transition to a discussion on applications, *i.e.*, actual applications and potential ideas for use-cases. Finally, we conclude the chapter with a reflection of the aforementioned topics in the form of a discussion.

### **7.2 Technical Challenges**

Like conventional FR, unconstrained faces *in the wild* [150] yield more difficult - imagery collected from sources outside a controlled laboratory environment is subject to more variations in pose, illumination, and scale. For faces, there are even more variables to further complicate the

problem, such as expression and age. Furthermore, preparing to run such benchmarks to mimic real-world use-cases (*i.e.*, designing experiments and preparing the data) is, in itself, a challenge. Inheriting these challenges, but adding even more variations inherent in nature and in true data distributions of kinship, it is unsurprising that visual kinship recognition is a difficult problem. Nonetheless, great efforts over the last decade have been spent not just on solving the problems in kinship recognition, but also critiquing kinship research and its direction. We now elaborate on the challenges to keep this technology from making the transition of research-to-reality.

### 7.2.1 Current limitations of SOTA

Still, we are close to achieving a performance-rating necessary for some applications (Section 7.3). From this, we perceive that bridging the gap between research-and-reality (*i.e.*, transitioning from research-to-practice) is happening. Upon a clear assessment of the state of progress in research, we highlight barriers still in need of overcoming, along with sharing edge cases as means of highlighting common errors. Hence, we aim to inspire by explicitly depicting weaknesses in current SOTA systems.

A clear limitation, however, is that most solutions for visual kinship recognition assume the relationship type *apriori*. Sometimes this could be practical, like if given a known source to decide whether or not the face, when paired with a target, is *KIN* or *NON-KIN*. Nonetheless, when considering the broader HCI incentive, along with data mining with social context, it is desirable to predict the exact type of relationship (*i.e.*, not just *KIN* or *NON-KIN*). Nonetheless, a high confidence in knowing whether a relationship does exist could serve as powerful prior knowledge when classifying the specific type.

Let us now consider the renowned KinFaceW dataset. Although the dataset has had a great impact in research, for having attracted many to the problem and, thus, has motivated many outstanding works, there are a few clear flaws in relating the results to real-world data. More than half of all true pairs making up KinFaceW are faces from the same photo. Researchers have then questioned the validity of the patterns being learned, showing that naive approaches such as color features [268] or detecting whether or not faces are from the same photo [12] outperform SOTA on most datasets, including KinFaceW. Thus, another clear limitation of some data resources is in the data distribution itself – a technical challenge we soon cover in-depth (Section 7.2.4).



**Figure 7.1: Sample faces synthesized to improve predictive power for faces of elderly adults (visualization from [11]).** Two models (*i.e.*, one per gender) were trained to synthesize input faces as younger– male fathers (*i.e.*, rows 1-2) and female mothers (*i.e.*, rows 3-4), the top sample is the original and the generated is below.

### 7.2.2 The nature

**Demographics and inherent bias.** A challenge are issues of bias in FR machinery. However, no study of bias in demographics (*e.g.*, ethnicity and gender) for kin-based data: a study that should be conducted, like Robinson *et al.* found variations in score sensitivities across subgroups in FR [25].

**Effects of age variations.** Family members with a large age-gap makes for more of a challenge. Wang *et al.* demonstrated a benefit in having a face image synthesized at younger ages [11]. Their ablation study revealed cumulative improvements as  $x \sim p_Y$  was bounded to  $>20$  years of age, then to  $>30$ , and up to  $>50$ . Improved results came with increasing the size of the domain (*i.e.*, the respective age considered young, which is orthogonal to those considered old). Figure 7.1 depicts samples of parents synthesized for kinship verification. Other augmentation techniques also proved useful, like transforming faces to their basis to then invert, rotate, and change ocular geometry [269].

### 7.2.3 The environment

The challenges from age variations in FR not only intensify in kin-based problems, but also change in novel ways. For instance, let us assume a comparison in the faces of a grandmother and a prospective grandson. The age of each and age gap between the two are subject to variation. In other words, the problem inherits the same challenges of FR such that considerations for directed relationships of concern– the grandmother might be in her early years when the picture was captured, just as the grandson might even be a grandfather himself at the time the picture was taken.



Figure 7.2: **FSP data (modified from [12]).** Data tagged via constraints, *must* and *cannot-link*:  $\approx 1M$  data points scraped from the web via 125 non-kin queries (*e.g.*, school student, sports team).

Nurture adds additional challenges to the problem: For instance, a pair of brothers inherited the nose from their mother; one boy experienced a broken nose perhaps more than once; suddenly, that boy no longer has a nose that resembles the mother. Where such challenges exist in conventional FR, the relative cost is greater with losing an inherited distinguishable feature from a prospective parent(s) in kin-based problems.

Biology-based research has focused on the problem of kinship recognition from a vast array of viewpoints. For instance, work that precedes the work done in machine vision, focused on a human's ability to recognize kinship— specifically, the ability of younger siblings to better distinguish between *KIN* and *NON-KIN* in strangers [270] (*i.e.*, having seen the first-born their entire lives trains them). An interesting hypothesis indeed, which is supported in the reported experiments (minimal sample set, but typical of human evaluations done in face-based research). Intuitively, the contrary could also be true (*i.e.*, the role of the older sibling, watching after their younger sibling would better train for this ability). In any case, the authors propose a theory conditioned on age; difference in age could play a significant role in such a study, as we agree this could be the case for a much older sibling (*i.e.*, already developing an ability to discriminate between faces), the same argument of realizing the key differences as a means of recognizing kin in a sibling at a young age could be argued both ways. Furthermore, the authors discarded samples of subjects with no siblings and more than two siblings— on the one hand the intent to control the experiment with less variation is understandable - on the other, subjects without siblings would serve as a meaningful baseline, while those with a number of siblings only strengthens the case for the oldest being the most keen on recognizing kin (*i.e.*, having grown watching over their younger siblings).

#### 7.2.4 The data and its distribution

Within-family variations are vast. As such, one cannot infer that the inherited traits from one father-son pair would mimic inherited traits of another father-son pair. Furthermore, the factors

introducing added complexity vary across different ethnic groups.

To capture the true data distributions of visual kinship as seen around the world is a great challenge, where many efforts have exhibited exploitable flaws. For instance, using color features claimed SOTA on the KinFaceW dataset, as faces of true-relatives often were cropped from the same photos [268][271]. The same motivation ushered in a different paradigm as means to measure unintended data leakage in the unnatural domain inherited by samples being of the same image or different. To say the least - this was a crafty piece of work that acquired an abundance of cheap data by image-level constraints that impose faces in the same photo as *matches*, which means it is a binary problem with classes for the *same* and *different* photo. In other words, by the paired data acquired by finding images with one-to-many faces from the web (Figure 7.2), Dawson *et al.* proposed training a detector to determine whether a face pair was from the *same* or *different* photo. Then, the boolean class model was directly evaluated on kin-based image sets, with the only difference in the target classes (*i.e.*, *same* and *different* assumed to be *KIN* and *NON-KIN*). Thus, showing SOTA ratings on a majority of existing kinship data– again, hypothesis that public benchmarks were subject to unintended data leakage, and one that is intrinsic to the distribution of classes (*i.e.*, *KIN* and *NON-KIN*). In the end, FSP proved competitive on KFW-I, KFW-II, Cornell KF, and TSKIN; however, FSP lacks sufficient training to perform well on the multi-image FIW data (*i.e.*, 58.6%, which was the first, smallest version of the FIW dataset). In fact, at the core of FIW specifications, as defined in its earliest paper [4], the concept of same and different photo was one considered in the creation of FIW– mentioned as part of motivation for the data in other recent literature reviews on kin-based image datasets [12].

### 7.3 Applications

We next review the use-cases for kinship recognition technology.

**Entertainment and personal knowledge.** AncestryDNA claimed >15 billion people in its DNA network: their >3M paying subscribers (and >16M people DNA tested), resulted in the establishment of 100M family trees that form 13B connections across 80 countries.<sup>1</sup> As of 2019, Ancestry launched AncestryHealth as a means to infer inheritable health conditions via DNA. Clearly, there is high interest in learning about one’s family roots– which started from curiosity (*i.e.*, knowing where one fits, recalling the aforementioned words of Furstenberg [103]), but now includes learning about one’s health from their DNA. Acquiring sufficient data to support DNA and imagery would be difficult.

---

<sup>1</sup>[www.ancestry.com/corporate/about-ancestry/company-facts](http://www.ancestry.com/corporate/about-ancestry/company-facts)

## CHAPTER 7. KINSHIP RECOGNITION - STATE OF TECHNOLOGY

However, provided more reliable kinship recognition capabilities, such technology would certainly enhance popular services such as those provided by billion dollar companies (*e.g.*, [ancestry.com](#)).

**Connect families.** Identify unknown children being exploited online; reconnect families separated by the modern-day refugee crisis [272]; find unknown relatives, whether directly or indirectly. Statistics show that people want to learn of missing family ties. Furthermore, unfortunate scenarios leave family members desperate to reconnect with lost member(s). Alternatively, law-enforcement could use kinship to solve other high-profile crimes— the decades long mystery of who the *Golden State Killer* was got solved by using DNA to build his family tree [273].

**Soft attribute as prior knowledge for traditional FR.** Whether it be to enhance FR capabilities [274], to learn to discriminate between hard negatives (*e.g.*, brothers), or to narrow the search (*e.g.*, FR failed to identify bombers of the 2013 Boston Marathon) - but had we known they were brothers, the search space could have been drastically reduced. Hence, kinship provides a powerful cue to help boost existing FR systems.

**Nature-based studies.** With the new millennium came the ability of 3D scans of facial appearances of ten pairs of twins to be compared via landmark features (*i.e.*, anteroposterior and vertical facial parameters) [275]. About ten years later, this inspired Dehghan *et al.* to ask: *Who do I look like?* And then attempt to solve the question using computer vision (*i.e.*, gated AE [8]).

**Kin-based face synthesis.** An early attempt to predict the appearance of a child from prospective parents was in [276]. Specifically, Froud *et al.* proposed EvoFit, which used classic shape-based modeling and *eigenfaces* to project a pair of faces via statistical appearance-based modeling. In all fairness, the generative task was heavily influenced by [57], as many face synthesis tasks were throughout the years, and especially in 2006 the EvoFit came out. In short, EvoFit learned its weights from face samples collected in a tightly controlled setting— per the requirement that 223 landmarks were precisely marked for all faces. As seminar as EvoFit was in its own right, this early attempt to predict the appearance of children was seemingly ahead of its time, in available machinery (lacking the data-driven, highly complex modeling techniques of today), in resources available to reproducible (*i.e.*, no public data released with paper), and in the problem statement itself. In other words, considering EvoFit was proposed before our 2010 timeline means it predated the first benchmark in kinship verification. With that, we believe the small impact of this work was due to its timing and, in return, the lack of complete support for the problem, so if others did want to partake they too would have to collect data. Meaning, it was impossible to reproduce results directly. Regardless how minimal the impact was in citations and usage of other researchers, the

work certainly showed promise considering the results were from a minimally-sized data pool. Thus, had a widely used benchmark been practiced, or provided the data constraints were handled (*i.e.*, inability to generalize + inability for others to reproduce), then EvoFit could have attracted much more attention. Perhaps, our 2010 time-line would have had to start a few years prior. Nonetheless, this is only speculation and, therefore, we can only hypothesize the *what ifs* after the fact.

## 7.4 Discussion

After we surveyed in [18], it was clear that a decade of research in visual kinship recognition resulted in an increased interest with an increase in data resources that were available. Clearly, the problems alone are challenging, even when compared to other machine-vision tasks (*e.g.*, conventional FR). Furthermore, the task of designing, collecting, and annotating labels is exceptionally difficult for kin-based problems. Thus, as contributions in data are proposed, interest seems to spike in response. With the release of the large-scale FIW dataset, for the first time, a data resource attempts to closely mimic data distributions of families around the globe. Moreover, FIW provides the data needed for the modern day, deep learning models. FIW, having had many existing datasets to learn from, remains the largest and most dynamic. However, the release of FIW was only the beginning, as efforts were then spent on annual challenges (*i.e.*, four consecutive years, 2017-2020, and also a Kaggle competition). With the resource and incentive provided by challenges, motivation for researchers to engage is ever so high and thus, we present this survey- not only as a means to realize the aspects that have been effective and vice versa- but also for a solid foundation for the next decade to build upon well-defined protocols and problem statements, each supported with source code, enabling even a wider audience to get started and contributing to the problem.

The deep learning revolution has only begun for visual kinship recognition - how to embed, how to fuse, how to interrupt - how do experts across disciplines engage by leveraging for a deeper understanding in inheritance from a strictly scientific point-of-view (*i.e.*, anthropology)! Hence, if we can devise the right tools for the right scholars synergy is bound to reveal insights in the nature of faces within families. Considering the many benchmarks that have a lot of room for improvement, along with the many social and relational data mining that is made possible with soft-attribute labels such as those in FIW, it is an exciting time for junior, senior, and practical researchers to reap benefits alongside its place with pure business, product, and patent design.

# Chapter 8

## Bias in Face Recognition

### 8.1 Overview

The more our society becomes integrated with machine learning (ML), the higher the interest in topics such as bias, fairness, and even the implications for the underlying formalization of existing or prospective ML standards [277, 278, 279]. Thus, an effect of vast companies growing more dependent on ML is an ever-increasing concern about the biased and unfair algorithms, *e.g.*, untrustworthy and prejudiced face recognition (FR) [280, 281].

A common trend in both research and mainstream has grown clear: the more we depend on technology that accelerates or automates everyday tasks, the more attention concepts such as biased and unfair algorithms should receive [282]. Furthermore, systems deployed for sensitive tasks, like biometrics [283] (*e.g.*, FR), need to be fully considered and understood. The perspective here recognizes the tendency of the researcher, the reporter, and the consumer to maintain transparency. Nowadays, in FR Convolutional Neural Networks (CNNs) are trained on a large number of faces identified by a detection system (Chapter 3). Recall that the goal is to encode faces in an N-dimensional space that pulls together samples of the same identity and pushes those that are different further apart. So, CNNs are trained to encode faces from the face (*i.e.*, image-space) to encoded representation (*i.e.*, feature-space). Face images are mapped to feature vectors and evaluated via a similarity scoring function, and the pairs with a similarity score above the decision threshold are assumed *genuine* and all others are classified as *imposters*.

Typically, a fixed threshold sets the decision boundary by which to compare scores (Figure 8.1). As such, features of the same identity must satisfy a criterion based on a single value [1, 48, 49, 172]. However, we found that an individual (*i.e.*, global) threshold– a crude measure

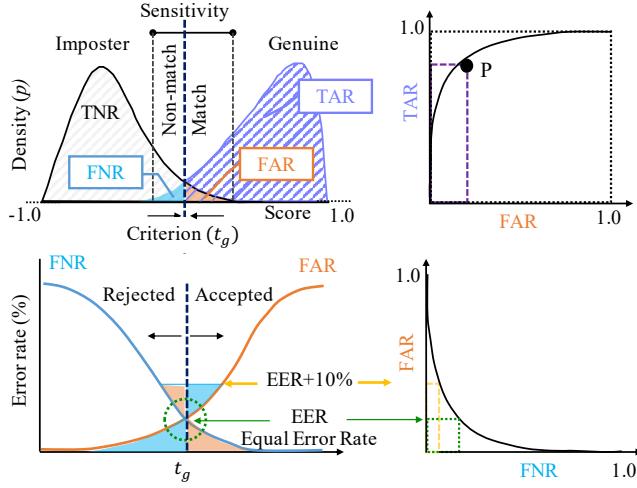


Figure 8.1: **Depiction of the biometrics.** The SDM shows the sensitivity related to a single threshold  $t_g$  (*top-left*). The area to the right of the threshold considers all accepted pairs, both correctly and incorrectly predicted. True Acceptance Rate (TAR) as a function of False Acceptance Rate (FAR) is a common way to report ratings for given false-rates (*top-right*). Equally common in FR is the trade-off between false-negative rate (FNR) and FAR (*bottom-left* and *bottom-right*).

that leads to skewed results across demographics and other attributes— are determined using a held-out set. In other words, the threshold for the FR *Matching* module is set according to the desired target: in research the threshold that yields the highest accuracy on the validation set; in practice the threshold is determined by the value that yields the desired rate, which depends on the use-case. For example, in the use-case where FR is used to enable entry via access control may have a smaller threshold that will realize fewer *genuine* samples than FR used for tagging photos per software recommendation—the falsely accepted instances of the latter use-case will have minimal, if any, negative effects. An important note to consider for the concept of having a held-out set is that it typically shares the same distribution with the test set, meaning it favors the same demographics as the held-out set had. That skew (*i.e.*, the difference in the performance of an algorithm of particular demographics) is our definition of bias. A key question is: *is FR too biased, or not?*

Now, provided two or more faces features encoded by a CNN, a distance (or similarity score)  $s$  must be learned such to act as a decision boundary to separate the genuine pairs score from the imposters score. Ideally, genuine and imposter scores would be completely separable. However, this is not the case in practice. It is this score-threshold (*i.e.*,  $\theta$ ) that determines whether or not the pair should be accepted. The implications are for faces features: to be assumed as the same, genuine

## CHAPTER 8. BIAS IN FACE RECOGNITION

Table 8.1: **Database stats and nomenclature.** *Header:* Subgroup definitions. *Top:* Statistics of Balanced Faces In the Wild (BFW). *Bottom:* Number of pairs for each partition. Columns grouped by ethnicity and then further split by gender.

	Asian (A)		Black (B)		Indian (I)		White (W)		
	Female (AF)	Male (AM)	BF	BM	IF	IM	WF	WM	Aggregated
# Faces	2,500	2,500	2,500	2,500	2,500	2,500	2,500	2,500	20,000
# Subjects	100	100	100	100	100	100	100	100	800
# Faces / Subject	25	25	25	25	25	25	25	25	25
# Positive Pairs	30,000	30,000	30,000	30,000	30,000	30,000	30,000	30,000	240,000
# Negative Pairs	85,135	85,232	85,016	85,141	85,287	85,152	85,223	85,193	681,379
# Pairs (Total)	115,135	115,232	115,016	115,141	115,287	115,152	115,223	115,193	921,379

class the score (or distance) must satisfy a criterion in the form of a single value [1, 49, 48, 172]. Mathematically, the decision  $D$  in similarity space is defined as

$$D = \begin{cases} \text{accept} & \text{if } s \geq \theta \\ \text{reject} & \text{if } s < \theta \end{cases}.$$

The importance of  $\theta$  should not be overlooked - a hyper-parameter that is a decision boundary in metric space. The optimal value depends on the specific use-case (*i.e.*, larger thresholds yield a lower probability that a sample is predicted as a true match). Regardless, the choice in threshold has a clear trade-off between false-positive (FP) and false-negative (FN) rates. For instance, a system that is claimed to perform at an error rate of 1 and 10,000, *i.e.*, one in every ten-thousand instances are incorrectly matched. We would then set our system by determining the threshold based on held-out data samples that allow the desired target error rate to be matched (Figure 8.1). The problem, per convention, is which assumes an average result across a held-out that then only holds true on the distribution of the source used. Then, specific cohorts (*e.g.*, ethnicity, gender, and other demographics) are unequally weighted due to an unequal representation. So, that single, global threshold, which is a sort of crude measure to begin with, is skewed to different cohorts that are not fairly represented by the source. In the end, these systems favor certain demographics, and it is the bias for which a change in cohort causes a change in the average performance of an algorithm.

Making matters more challenging, precise definition of race and ethnicity vary from source-

## CHAPTER 8. BIAS IN FACE RECOGNITION

to-source. For example, the US Census Bureau allows an individual to self-identify race.<sup>1</sup> Even gender, our attempt to encapsulate the complexities of the sex of a human as one of two labels. Others have addressed the oversimplified class labels by representing gender as a continuous value between 0 and 1 - rarely is a person entirely  $M$  or  $F$ , but most are somewhere in between [221]. For this work, we define subgroups as specific sub-populations with face characteristics similar to others in a region. Specifically, we focus on 8 subgroups (Figure 8.2).

The adverse effects of a global threshold are two-fold: (1) mappings produced by CNNs are nonuniform. Therefore, distances between pairs of faces in different demographics vary in distribution of similarity scores (Fig 8.3); (2) evaluation set is imbalanced. Subgroups that make up a majority of the population will carry most weight on the reported performance ratings. Reported results favor the common traits over the underrepresented. Demographics like gender, ethnicity, race, and age are underrepresented in most public datasets [221, 279].

For (1), we propose subgroup-specific (*i.e.*, optimal) thresholds while addressing (2) with a new benchmark dataset to measure bias in FR, BFW (Table 8.2 and 8.4). BFW serves as a proxy for fair evaluations for FR while enabling per subgroup ratings to be reported. We use BFW to gain an understanding of the extent to which bias is present in state-of-the-art (SOTA) CNNs used FR. Then, we suggest a mechanism to mitigate problems of bias with more balanced performance ratings for different demographics. Specifically, we propose using an adaptive threshold that varies depending on the characteristics of detected facial attributes (*i.e.*, gender and ethnicity). We show an increase in accuracy with a balanced performance for different subgroups. Similarly, we show a positive effect of adjusting the similarity threshold based on the facial features of matched faces. Thus, selective use of similarity thresholds in current SOTA FR systems provides more intuition in FR research with a method easy to adopt in practice.

The contributions of this work are three-fold. (1) We built a balanced dataset as a proxy to measure verification performance per subgroup for studies of bias in FR. (2) We revealed an unwanted bias in scores of face pairs - a bias that causes ratings to skew across demographics. For this, we showed that an adaptive threshold per subgroup balances performance (*i.e.*, the typical use of a global threshold unfavorable, which we address via optimal thresholds). (3) We surveyed humans to demonstrate bias in human perception (NIH-certified, *Protect Humans in Research*).

The adverse effects of a global threshold are three-fold: (1) the evaluation set is typically imbalanced. The demographics of the majority are weighted more in the reported performance ratings. Therefore, reported results skew to rarer traits that are more common in the underrepresented

---

<sup>1</sup> [www.census.gov/mso/www/training/pdf/race-ethnicity-onepager](http://www.census.gov/mso/www/training/pdf/race-ethnicity-onepager)

## CHAPTER 8. BIAS IN FACE RECOGNITION

subgroups— a phenomena that should be considered for different subgroups (*i.e.*, gender, ethnicity, race, age). (2) the mappings produced by a CNN have various levels of sensitivity in the metric (Figure 8.3). Therefore, the range of distances between true pairs varies across demographics. (3) a global threshold is referenced when comparing face encodings. Since the optimal score shifts for different demographics, there ought to be variable thresholds (*i.e.*, sliding threshold set according to demographic information). Furthermore, the validation set used to determine global threshold and the test set used to report results should be understood - this issue will be unnoticed in performance ratings if the validation and test are from the same distribution and, thus, the resulting performance ratings are based in favor of the majority. This leads to performance ratings that are incomplete and even misleading.

To address the lack of a balanced data (*i.e.*, (1)), we propose to evaluate on our dataset, which was built specifically for measuring biases in demographics for facial verification (FV) systems in a systematic, reproducible way. With it, we introduce a new benchmark for FR, called BFW (Figure 8.7). BFW serves as a platform to fairly evaluate FR systems and enable demographic-specific ratings to be reported (Table 8.4). We use BFW to gain a deeper understanding of the extent of bias present in facial embeddings extracted from a SOTA CNN model. We then suggest a mechanism to counter the biased feature space to mitigate problems of bias with more balanced performance ratings across demographics, and all the while improving the overall accuracy. Specifically, we unlearn demographic knowledge in face encodings, while preserving identity information. Thus, we learn to map the encodings to a lower dimensional space containing less knowledge of subgroups. The byproducts are then embeddings that preserve the privacy of its subject’s ethnicity and gender (*i.e.*, subgroups). It is this feature adaptation scheme proposed to address items (2-3).

Our contributions in topics of bias (*i.e.*, this [25]) are the following:

- Demonstrate a bias in an existing SOTA CNN with our BFW dataset (Table 8.2). We propose a feature learning scheme that employs domain adaptation to debias face encodings and, most importantly, balances performances across subgroups such to boost the overall performance.
- Hide attribute information in encodings— a byproduct of the proposed debiasing scheme is the reducing knowledge of attributes. Beyond privacy, it removes other potential biases, whether unintended (*e.g.*, models trained on top) or intended bias (*e.g.*, human consciously using).
- Provide insights with analysis of hard samples overcome by the proposed debiasing scheme. Evidence in the form of salience mapping and face pairs are shown and discussed.

## CHAPTER 8. BIAS IN FACE RECOGNITION

- Develop code-base as public Git Hub (*i.e.*, <https://github.com/visionjo/facerec-bias-bfw>); provided form (*i.e.*, link <https://forms.gle/3HDBikmz36i9DnFf7>) for dataset download requests, where paired data and related resources used in *facerec-bias-bfw* repo are available.

### 8.1.1 Organization

The rest of the chapter is organized as follows. In Section 8.2, we review work related to bias in FR, along with works related to the problem and solution spaces. We then cover our BFW dataset— the motivation, specifications, and described (Section 8.3). Then, in Section 8.4, we introduce the proposed methodology. Section 8.6 follows this with settings and results of the experiments, along with the details of our BFW database. Finally, we conclude and discuss next steps in Section 8.7.

## 8.2 Related Work

We next review the research related to bias and privacy in FR, both for humans and machines, and along with some background information required to understand the motivation and overall solution of the proposed model. Specifically, we first briefly discuss problems of bias in general ML, then that which is specific to FR. Following this, we support our hypothesis that human too possess a similar bias (*i.e.*, more familiar to those of subgroup most frequently seen in the past). Then, we describe several works in domain adaptation— the domain for which our proposed solution is best characterized. Finally, we cover problems of privacy in FR.

### 8.2.1 Bias in machine learning

The progress and commercial value of ML are exciting. However, due to inherent biases in ML, society is not readily able to trust completely in its widespread use. The exact definitions and implications of bias vary between sources, as do its causes and types. A common theme is that bias hinders performance ratings in ways that skew to a particular sub-population. In essence, the source varies, whether from humans [284], data or label types [285], ML models [286, 287], or evaluation protocols [288]. For instance, a vehicle-detection model might miss cars if training data were mostly trucks. In practice, many ML systems learn biased data, which could be detrimental to society.

### 8.2.2 Bias in facial recognition

Biases in FR focus on characterizing performance across various *soft attributes*, such as gender, ethnicity, or age [283]. Researchers have spent great efforts proposing problem statements and solutions to problems of bias in FR technology. We focus on the demographics (or subgroups) of gender and ethnicity. The inherent problems here are two-fold. First, gender is handled as a boolean label, which is a gross approximation of individual uniqueness in question of sexuality - a spectrum of real numbers would be more appropriate [221]. Secondly, the definitions of race and ethnicity are loosely defined. The US Census Bureau allows an individual to self-identify race.<sup>2</sup> We define it as a group of people having facial characteristics similar to those found in a region. The result is various types of biases in FR systems in favor of or against particular demographics remain a question.

Balakrishnan *et al.* trained a generator to manipulate latent space features in that the controlled attributes were skin-tone, length of hair, and hair color[289]. Another synthesis solution proposed was to generate faces across various ages as a means to augment training data [290]. Terhorst *et al.* recognized the same phenomena our work is found on– the variation in sensitivities of scores for different demographics [291]. Specifically, the authors propose a score normalization scheme to handle the problem of inaccurate performance ratings when demographic-specific performances are compared to the average– a problem highlighted in paper [25].

Some aim to characterize the amount of bias in a system, whether it be for gender [292, 293, 294], ethnicity, age [295], or two or more of the aforementioned [280, 296, 297, 298, 299]. A recent *European Conference on Computer Vision* (ECCV) challenge provided incentive and for researchers to propose solutions for problems of bias with respect to ethnicity, gender, age, pose, and with and without sunglasses [300]. Other recent works in FR technology introduce additional modalities, such as profile information, to the problem of bias [301, 302]. Another research question concerns the measuring of biases in FR systems [296, 292]. Some focus on templates [303]. Some debias at the score level [291]. Some focus to debias pre-trained models [304]. Wang *et al.* introduced a reinforcement learning-based race balance network (RL-RBN) to find optimal margins for non-Caucasians as a Markov decision process before passing to the deep model that learns policies for an agent to select margins as an approximation of the Q-value function; *i.e.*, the skewness of feature scatter between races can be reduced [305]. Even more HCI-based views have been introduced as semi-supervised bias detection systems that act as tools with humans in-the-loop [306].

Yin *et al.* proposed to augment the feature space of underrepresented classes using

---

<sup>2</sup><https://www.census.gov>

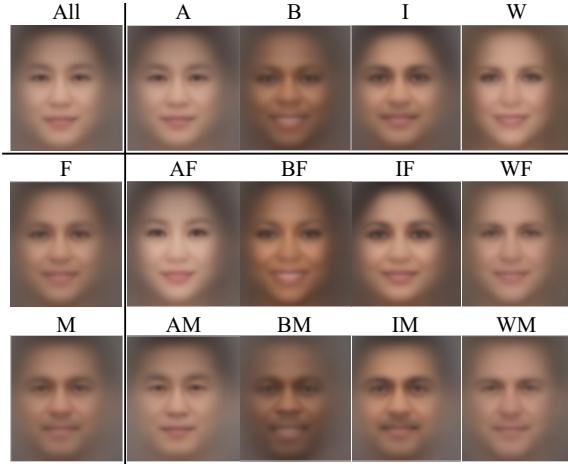


Figure 8.2: **BFW dataset.** Average face per subgroup: *top-left*: the entire BFW; *top-row* per ethnicity; *left-column*: per gender. The others represent the ethnicity and gender, respectively. Table 8.1 defines the acronyms of subgroups.

different classes with a diverse collection of samples [307]. This was to encourage distributions of underrepresented classes to resemble the others more closely. Similarly, others formulated the imbalanced class problem as one-shot learning, where a generative adversarial network (GAN) was trained to generate face features to augment classes with fewer samples [87]. Generative Adversarial Privacy and Fairness (GAPF) was proposed to create fair representations of the data in a quantifiable way, allowing for the finding of a de-correlation scheme from the data without access to its statistics [308]. Wang *et al.* defined subgroups at a finer level (*i.e.*, Chinese, Japanese, Korean), and determined the familiarity of faces inter-subgroup [309]. Genders have also been used to make subgroups (*e.g.*, for analysis of gender-based face encodings [310]). Most recently, [279] proposed to adapt domains to bridge the bias gap by knowledge transfer, which was supported by a novel data collection, Racial Faces in-the-Wild: (RFW). The release of RFW occurred after BFW was built - although similar in terms of demographics, RFW uses faces from MSCeleb [196] for testing, and assumes CASIA-Face [149] and VGG2 [14] were used to train. In contrast, our BFW assumes VGG2 as the test set. Furthermore, BFW balance subgroups: RFW splits subgroups by gender and race, while BFW has gender, race, or both).

Most similar to us is [295, 311, 312, 313] - each was motivated by insufficient paired data for studying bias in FR. Then, problems were addressed using labeled data from existing image collections. Uniquely, Hupont *et al.* curated a set of faces based on racial demographics (*i.e.*, *Asian* (A), *Black* (B), and *White* (W)) called Demographic Pairs (DemogPairs) [312]. In contrast, [295]

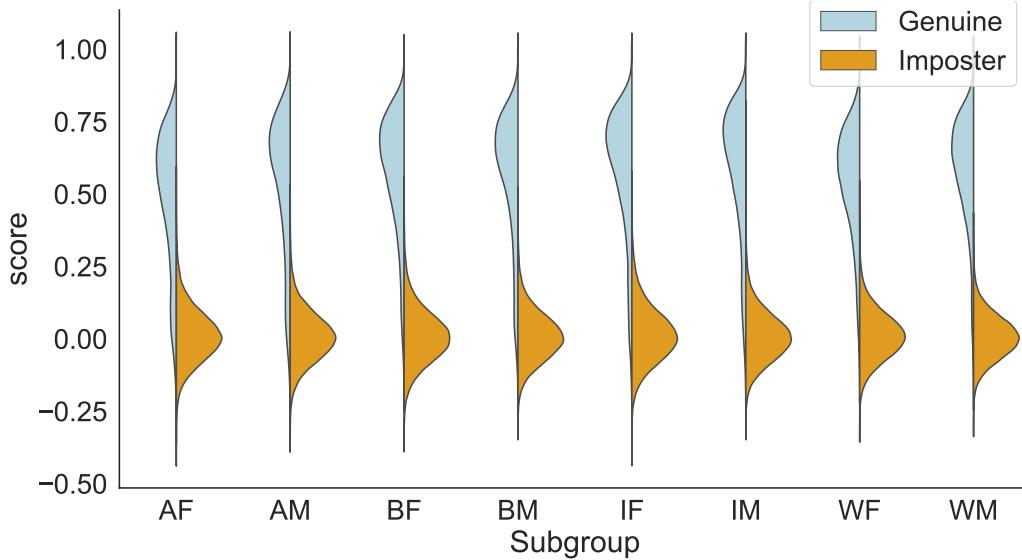


Figure 8.3: **Signal detection model (SDM) across subgroups.** Scores of *imposters* have medians  $\approx 0.3$  but with variations in upper percentiles; *genuine* pairs vary in mean and spread (*e.g.*, AF has more area of overlap). A threshold varying across different subgroups yields a constant FAR.

honed in on adults versus children called Wild Child Celebrity (ITWCC). Like the proposed BFW, both were built by sampling existing databases, but with the addition of tags for the respective subgroups of interest. Aside from the additional data of BFW (*i.e.*, added subgroup *Indian* (I), along with other subjects with more faces for all subgroups), we also further split subgroups by gender. Furthermore, we focus on the problem of facial verification and the different levels of sensitivity in cosine similarity scores per subgroup.

### 8.2.3 Human bias in machine learning

Bias is not unique to ML - humans are also susceptible to a perceived bias. Biases exist across race, gender, and even age [277, 314, 315, 316]. Wang *et al.* showed machines surpass human performance in discriminating between Japanese, Chinese, or Korean faces by nearly 150% [309], as humans just pass random (*i.e.*, 38.89% accuracy).

We expect the human bias to skew to their genders and races. For this, we measure human perception with face pairs of different subgroups (Section 8.3.3). The results concur with [309], as we also recorded overall averages below random (<50%).

### 8.2.4 Imbalanced data and data problems in FR

The impact from the quality of fairness depends on the context for which it is used. Furthermore, various paradigms have been proposed as means to a solution: some alter the data distribution to yield classifiers of equal performances for all classes (*i.e.*, re-sampling, like by under-sampling and over-sampling [317]); others alter the data itself (*i.e.*, algorithms that adjust classification costs). For instance, Oquab *et al.* re-sampled at the image patch-level [318]. Specifically, the aim was to balance foregrounds and backgrounds for object recognition. On the other hand, Rudd *et al.* proposed the mixed objective optimization network (MOON) architecture [319] that learns to classify attributes of faces by treating the problem as a multi-task (*i.e.*, a task per attribute) attribute to more balanced performances when training on data that has an imbalanced distribution across attributes. The Cluster-based Large Margin Local Embedding (CLMLE) [320] sampled clusters of samples in the feature-space that were used to regularize the models at the decision boundaries of underrepresented classes. See literature reviews for details on approaches that alter at the data or algorithmic-level [321, 322, 323].

More specific to faces, Drozdowski *et al.* summarizes that the cohorts of concern in biometrics are demographic (*e.g.*, sex, age, and race), person-specific (*e.g.*, pose or expression [324], and accessories like eye-wear or makeup), and environmental (*e.g.*, camera-model, sensor size, illumination, occlusion) [283]. Albiero *et al.* found empirical support that having training data that is well balanced in gender does not mean that results of a gender-balanced test set will be balanced [325]. Our studies focus on the effect demographic has in FV by assessing demographic-specific classification ratings. Our BFW data resource allows us to analyze existing SOTA deep CNNs on different demographics (or subgroups). We provide practical insight: FR benchmarks often report with misleading ratings—ratings are dependent on the demographics of the population.

To match the capacity of modern-day deep models several large FR datasets were released [196, 142, 14, 242]. More recently, several have reported on the imbalance in demographics are the data, and proposed balanced resources for FR-based tasks [326, 279, 327, 221]. Diversity in Faces (DiF) came first [221], which did not include identity labels. Moreover, DiF is no longer available for download. Others released data with demographics balanced, but for the task of predicting the demographic and, thus, do not include identity labels [279, 327]. Hupont *et al.* proposed DemogPairs is balanced across 6 subgroups of 600 identities from CASIA-WebFace (CASIA-W) [149], VGG [142], and VGG2 [14]. Similar to DemogPairs, except with 8 subgroups, 800 identities, and with number of faces per identity the same for all, our BFW data was the latest release for measuring

bias in FR technology. Furthermore, recognizing that existing SOTA models are already trained on a public resource, we built BFW from just VGG2 to minimize conflicts in overlap between train and test. Table 8.4 compares our data with the others.

### 8.2.5 Domain adaptation and feature alignment

Domain adaptation (DA) [328, 329, 330] employs labeled data from the source domain to make it generalize well to the typically label-scarce target domain; hence, a common solution to relieve annotation costs. DA can be roughly classified as the semi-supervised DA [330] or the unsupervised one [331] according to the access to target labels. The crucial challenge toward DA is the distribution shift of features across domains (*i.e.*, domain gap), which violates the distribution-sharing assumption of conventional machine learning problems. In our case, the different domains are the different subgroups (Table 8.2).

To bridge such gap, some of feature alignment (FA) methods attempt to project the raw data into a shared subspace where certain feature divergence or distance is minimized to confuse them. Various methods, such as Correlation Alignment (CORAL) [332], Maximum Mean Discrepancy (MMD) [333], and Geodesic Flow Kernel (GFK) [334, 335], have been developed in this line. Currently, adversarial domain alignment methods (*i.e.*, DANN [336], ADDA [337]) have attracted increasing attention by designing a zero-sum game between a domain classifier (*i.e.*, discriminator) and a feature generator. The features of different domains will be mixed if the discriminator cannot differentiate the source and target features. More recently, learning well-clustered target features proved to be helpful in conditional distributions alignment. Both DIRT-T [331] and MME [330] methods applied entropy loss on target features to implicitly group them as multiple clusters in the feature space to keep the discriminative structures through adaptation. Inspired by FA, we aim to align the score distributions of different subgroups by adjusting score sensitivities (Figure 8.3).

### 8.2.6 Protecting demographic information in FR

ML is growing more accessible. As such, it grows more in our day-to-day lives. The levels of sensitivities of the use-case ought to be put under careful consideration - not only the model, but the data too [338]. That means, intermediate results, also known as features are included. FR is an ML problem that has made great progress in recent years, being used *off-the-shelf* in many applications. It is time to carefully consider data privacy concerns - with priority on topics of biometrics. Provided

careless or too little action is spent to protect the user behind the face image, the more the chance that data may be used by an adversary maliciously [339].

Several have recently attempted to solve problems of bias in demographic-based classifiers (*e.g.*, ethnicity and gender classifiers). Furthermore, attempts to disguise demographic information in facial encodings while preserving FV abilities have been proposed for privacy and protection purposes [340, 341]. In other words, prior works recognized the importance of preserving the identity information in facial features, while ridding it of evidence of demographics. Our model inherently does this as part of the objective aims for an inability to recognize subgroups.

The aforementioned assume the facial encoding are accessible - this makes sense in terms of reduced computations (*i.e.*, no need to encode each time) and storage (*i.e.*, encodings are smaller representations of the image). However, several works aimed to hide attribute information in image space; for instance, Othman *et al.* learned to morph faces to suppress gender information in the image-space while preserving the identification [342]. Guo *et al.* proposed a mapping from image-to-noise, both encrypting the image such that the encoder still decodes the identity but without the ability to determine gender by machine or human [343]. Ma *et al.* viewed the communication between servers as the point of concern for privacy- a lightweight privacy preserving adaptive boosting (AdaBoost) FR framework (POR) based on additive fusion for secret sharing to encrypt model parameters, while using a cascade of classifiers to address different protocols [344].

### 8.3 Balanced Faces In the Wild (BFW)

BFW provides balanced data across ethnicity (*i.e.*, Asian (A), Black (B), Indian (I), and White (W)) and gender (*i.e.*, Female (F) and Male (M)) – a total of eight demographics referred to as subgroups (Figure 8.7). As listed in Table 8.4, BFW has an equal number of subjects per subgroup (*i.e.*, 100 subjects per subgroup) and faces per subject (*i.e.*, 25 faces per subject). Note that the key difference between BFW and DemogPairs is the additional attributes and the increase in labeled data; the differences from RFW and FairFace are in the identity labels and distributions (Table 8.2).

BFW was built with VGG2 [14] by using a set of classifiers on the list of names, and then the corresponding face data. Specifically, we ran a name-ethnicity classifier [345] to generate the initial list of subject proposals. Then, the list was further refined by processing the corresponding faces with ethnicity [346] and gender [347] classifiers. Next, we manually validated, keeping only those that were true members of the respective subgroup. Faces for each subject were then limited to a total of 25 faces that were selected at random, with the distribution of the resolution of the detected

faces (*i.e.*, area of the bounding boxes) shown in Figure 8.6. Thus, BFW was obtained with minimal human input, having had the proposal lists generated by automatic machinery.

The subgroups of BFW were determined based on physical features most common amongst the respective subgroup [25], which can be regarded as multiple domains due to the feature distributions mismatch across these subgroups. However, the assumption that a discrete label has the capacity to describe an individual is, at best, imprecise. Nonetheless, the assumption allows for a finer-grain analysis of subgroup and is a step in the right direction. Thus, we refute any claim that our efforts here are the final solution; however, we insist that the data and proposed machinery are merely an attempt to establish a foundation for future work to extend. In any case, the two gender for the four ethnic groups make up the eight subgroups of the BFW dataset (Figure 8.7). Formally put, the tasks addressed have labels for gender  $l^g \in \{F, M\}$  and ethnicity  $l^e \in \{A, B, I, W\}$ , where the  $K$  subgroups (*i.e.*, demographics) are then  $K = |l_g| * |l_e| = 8$ .

We next review the details behind the process of building BFW without any financial burden, while maintaining limited amounts of human labor requirements.

### 8.3.1 The data

Problems of bias in FR motivated us to build BFW. Inspired by DemogPairs [312], the data is made up of evenly split subgroups, but with an increase in subgroups (*i.e.*, IF and IM), subjects per subgroup, and face pairs (Table 8.1 and Figure 8.2).

**Compiling subject list.** Subjects were sampled from VGG2 [14] - unlike others built from multiple sources, BFW has fewer potential conflicts in train and test overlap with existing models. To find candidates for the different subgroups, we first parsed the list of names using a pre-trained ethnicity model [345]. This was then further refined by processing faces using ethnicity [346] and gender [347] classifiers. This resulted in hundreds of candidates per subgroup, which allowed us to manually filter 100 subjects per the 8 subgroups.

**Detecting faces.** Faces were detected using MTCNN [102].<sup>3</sup> Then, faces were assigned to one of two sets. Faces within detected bounding box (BB) regions extended out 130% in each direction, with zero-padding as the boundary condition made-up one set. The second set were faces aligned and cropped for Sphereface [48] (see the next step). Also, coordinates of the BB and the five landmarks from *multi-task CNN* (MTCNN) were stored as part of the static, raw data. For samples with multiple

---

<sup>3</sup><https://github.com/polarisZhao/mtcnn-pytorch>

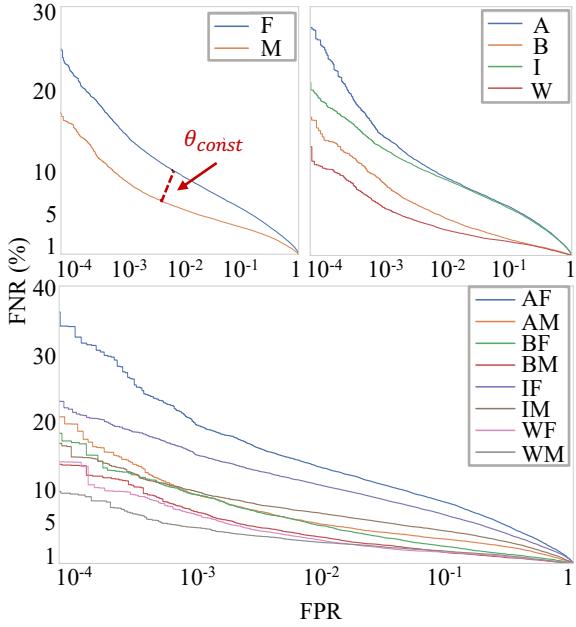


Figure 8.4: **Detection error trade-off (DET) curves.** *Top-left:* per gender. *Top-right:* per ethnicity. *Bottom:* per subgroup (*i.e.*, combined). Dashed line shows about  $2\times$  difference in FPR for the same threshold  $\theta_{const}$ . FNR is the match error count (closer to the bottom is better).

face detections, we used the BB area times the confidence score of the MTCNN to determine the face most likely to be the subject of interest, with the others set aside and labeled *miss-detection*.

**Validating labels.** Faces of BFW were encoded using the original implementation of the SOTA Sphereface [48]. For this, each face was aligned to predefined eye locations via an affine transformation. Then, faces were fed through the CNN twice (*i.e.*, the original and horizontally flipped), with two features fused by average pooling (*i.e.*, 512 D). A matrix of cosine similarity scores was then generated for each subject and removed samples (*i.e.*, rows) with median scores below threshold  $\theta = 0.2$  (set manually). Mathematically, the  $n^{th}$  sample for the  $j^{th}$  subject with  $N_j$  faces was removed if the ordinal rank of its score  $n = \frac{P \times N}{100} \geq \theta$ , where  $P = 50$ . In other words, the median (*i.e.*, 50 percentile) of all scores for a face with respect to all of faces for the respective subject must pass a threshold of  $\theta = 0.2$ ; else, the face is dropped. This allowed us to quickly prune FP face detections. Following [4, 6], we built a JAVA tool to visually validate the remaining faces. For this, the faces were ordered by decreasing confidence, with confidence set as the average score, and then displayed as image icons on top toggling buttons arranged as a grid in a sliding pane window. Labeling then consisted of going subject-by-subject and flagging faces of *imposters*.

**Sampling faces and creating folds.** We created lists of pairs in five-folds with subjects split evenly per person and without overlap across folds. Furthermore, a balance in the number of faces per subgroup was obtained by sampling twenty-five faces at random from each. Next, we generated a list of all the face pairs per subject, resulting in  $\sum_{l=1}^L \sum_{k=1}^{K_l} \binom{N_k}{2}$  positive pairs, where the number of faces of all  $K_l$  subjects  $N_k = 25$  for each of the  $L$  subgroups (Table 8.1). Next, we assigned subjects to a fold. To preserve balance across folds, we sorted subjects by the number of pairs and then started assigning to alternating folds from the one with the most samples. Note, this left no overlap in identity between folds. Later, a negative set from samples within the same subgroup randomly matched until the count met that of the positive. Finally, we doubled the number with negative pairs from across subgroups but in the same fold.

### 8.3.2 Problem formulation

FV is the special case of the two-class (*i.e.*, boolean) classification. Hence, pairs are labeled as the “same” or “different” *genuine* pairs (*i.e.*, *match*) or *impostor* (*i.e.*, *mismatch*), respectively. This formulation (*i.e.*, FV) is highly practical for applications like access control, re-identification, and surveillance. Typically, training a separate model for each unique subject is unfeasible. Firstly, the computational costs compound as the number of subjects increase. Secondly, such a scheme would require model retraining each time a new person is added. Instead, we train models to encode facial images in a feature space that captures the uniqueness of a face, to then determine the outcome based on the output of a scoring (or distance) function. Formally put:

$$f_{\text{boolean}}(\vec{x}_i, \vec{x}_j) = d(\vec{x}_i, \vec{x}_j) \leq \theta, \quad (8.1)$$

where  $f_{\text{boolean}}$  is the *matcher* of the feature vector  $\vec{x}$  for the  $i^{\text{th}}$  and  $j^{\text{th}}$  sample [150].

Cosine similarity is used as the *matcher* in Eq 8.1 the closeness of  $i^{\text{th}}$  and  $j^{\text{th}}$  features, *i.e.*,  $s_l = \frac{f_i \cdot f_j}{\|f_i\|_2 \|f_j\|_2}$  is the closeness of the  $l^{\text{th}}$  pair.

### 8.3.3 Human assessment

To focus on the human evaluation experiment, we honed-in on pairs from two groups, White Americans (W) and Chinese from China (C). This minimized variability compared to the broader groups of Whites and Asians, which was thought to be best, provided only a small subset of the data was used on fewer humans than subjects in BFW.

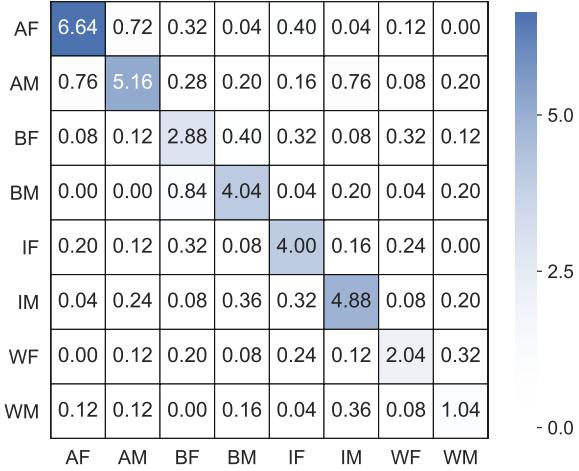


Figure 8.5: **Confusion matrix.** Error (Rank 1, %) for all BFW faces versus all others. Errors concentrate intra-subgroup - consistent with the SDM (Figure 8.3). Although subgroups are challenging to define, this shows the ones chosen are meaningful for FR.

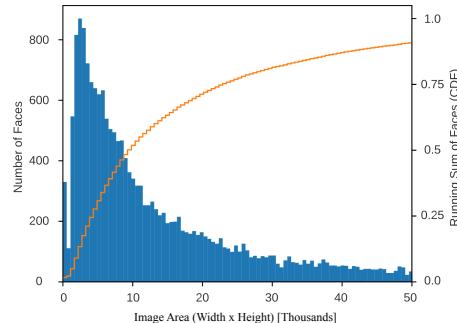
Samples were collected from individuals recruited from multiple sources (*e.g.*, social media, email lists, family, and friends) - a total of 120 participants were sampled at random from all submissions that were completed and done by a W or C participant. Specifically, there were 60 W and 60 C, both with *Male* (M) and *Female* (F) split evenly. A total of 50 face pairs of non-famous “look-alikes” were collected from the internet, with 20 (WA) and 20 (C) pairs with, again, M and F split evenly. The other 10 were of a different demographic (*e.g.*, Hispanic/ Latino, Japanese, African). The survey was created, distributed, and recorded via [PaperForm](#). It is important to note that participants were only aware of the task (*i.e.*, to verify whether or not a face-pair was a *match* or *non-match*, but with no knowledge of it being a study on the bias).

## 8.4 Methodology

To discuss the bias and privacy concerns addressed by the proposed, we first introduce facial verification (FV). Specifically, we review the problem statement, the supporting facial image dataset, along with the objectives of the proposed framework set to achieve the solutions sought in this work. That is to say, to preserve identity information while balancing the sensitivities of encodings for the different demographics (*i.e.*, subgroups), and in a way, to remove knowledge of the subgroups from the facial encoding - the typical representation available for operational cases.

**Table 8.2: Data statistics, notation, and scores for subgroups of our BFW data.** *Top:* Specifications of BFW and subgroup definitions. *Middle:* Number of pairs. *Bottom:* Accuracy fo a global threshold  $t_g$ , the value of the optimal threshold  $t_o$ , and accuracy using  $t_o$  per subgroup. Columns grouped by race and then further split by gender. Notice the inconsistent ratings across subgroups.

	Asian (A)		Black (B)		Indian (I)		White (W)		
	Female (AF)	Male (AM)	BF	BM	IF	IM	WF	WM	Aggregated
No. Faces	2,500	2,500	2,500	2,500	2,500	2,500	2,500	2,500	20,000
No. Subjects	100	100	100	100	100	100	100	100	800
No. Faces / subject	25	25	25	25	25	25	25	25	25
No. Positive	30,000	30,000	30,000	30,000	30,000	30,000	30,000	30,000	240,000
No. Negative	85,135	85,232	85,016	85,141	85,287	85,152	85,223	85,193	681,379
Total	115,135	115,232	115,016	115,141	115,287	115,152	115,223	115,193	921,379
Acc@ $t_g$	0.876	0.944	0.934	0.942	0.922	0.949	0.916	0.918	0.925±0.022
$t_o$	0.235	0.274	0.267	0.254	0.299	0.295	0.242	0.222	0.261±0.025
Acc@ $t_o$	0.916	0.964	0.955	0.971	0.933	0.958	0.969	0.973	0.955 ± 0.018



**Figure 8.6: BFW statistics (i.e., pixel counts).** Histogram of image areas in pixels (blue plot). The orange curve shows the cumulative count of images up to a given area.

### 8.4.1 Problem statement

FV systems make decisions based on the likeliness a pair of faces are of the same identity. In fact, the core procedure of verification systems are often similar to the FR employed for various applications. Specifically, a CNN is trained on a closed set of identities, and then later used to encode faces (*i.e.*, map face images to features). The encodings are then compared in closeness to produce a single score— often via cosine similarity in FR [177]. It is imperative to learn the optimal score

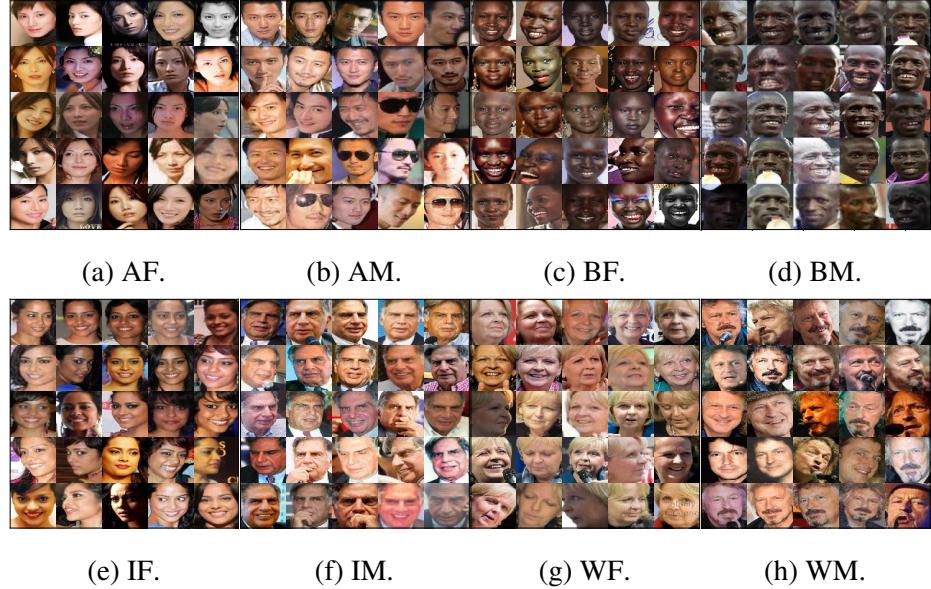


Figure 8.7: **Samples of BFW.** Per subgroup: the 25 samples for a random subject are shown.

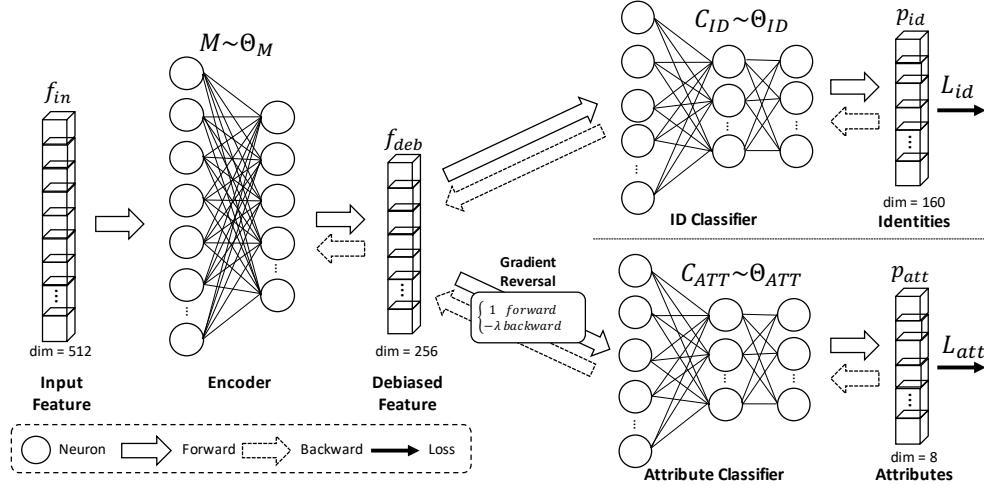
threshold separating true from false pairs. The threshold is the decision boundary in score space, *i.e.*, the *matching function*.

#### 8.4.1.1 The matching function

A real-valued similarity score  $R$  assumes a discrete label of  $Y = 1$  for *genuine* pairs (*i.e.*, a true-match) or  $Y = 0$  for *impostor* (*i.e.*, untrue-match). The real-value is mapped to a discrete label by  $\hat{Y} = \mathbb{I}\{R > \theta\}$  for some predefined threshold  $\theta$ . The aforementioned can be expressed as *matcher*  $d$  operating as

$$f_{\text{boolean}}(\vec{x}_i, \vec{x}_j) = d(\vec{x}_i, \vec{x}_j) \leq \theta, \quad (8.2)$$

where the face encodings in  $\vec{x}$  being the  $i^{\text{th}}$  and  $j^{\text{th}}$  sample – a conventional scheme in the FR research communities [150]. We use cosine similarity as the *matcher* in Eq. 8.2, which produces a score in closeness for the  $i^{\text{th}}$  and  $j^{\text{th}}$  faces (*i.e.*,  $l$ -th face pair) by  $d(\vec{x}_i, \vec{x}_j) = s_l = \frac{f_i \cdot f_j}{\|f_i\|_2 \|f_j\|_2}$ . The decision boundary formed by threshold  $\theta$  controls the level of *acceptance* and *rejection*. Thus,  $\theta$  inherits a trade-off between the sensitivity and specificity. The operating point chosen is done so, most always, depending on the purpose of the respective system (*i.e.*, perhaps higher sensitivity (or lower threshold) for threat-related tasks, in that flagging a few extra is worth not overlooking the one true



**Figure 8.8: Debiasing framework.** The framework used to learn a projection that casts facial encodings to a space that (1) preserves identity information (*i.e.*,  $C_{ID}$ ) and (2) removes knowledge of subgroup (*i.e.*,  $C_{ATT}$ ). The benefits of this are two-fold: ability to verify pairs of faces fairly across attributes and an inability to classify attribute for privacy and safety purposes. Note, that the *gradient reversal* [13] flips the sign of the error back-propagated from  $C_{ATT}$  to  $M$  by scalar  $\lambda$  during training.

positive). Specifically, the trade-off involves FNR, a Type 1 Error where *genuine* attempts to pass but is falsely rejected. Mathematically, it relates by

$$FNR = \frac{FN}{P} = \frac{FN}{FN + TP} = 1 - TAR = 1 - \frac{TP}{FN + TP},$$

with positive counts P, TP, and FN.

The other error type contributes to the FPR, is referred to as the Type II Error as an *impostor*, and is falsely accepted:

$$FAR = \frac{FP}{N} = \frac{FP}{FP + TN} = 1 - TNR = 1 - \frac{TN}{FP + TN},$$

where the counts negatives is N, and the metrics are TN, FP, TNR, and FAR. A *matching* function is the module typically characterized using the listed metrics (Figure 8.1). The geometric relationships of the metrics related to the score distributions and the choice in threshold show the trade-offs in error rates (Figure 8.1).

The hyperparameter  $\theta$  is often determined for a desired error rate on a held-out validation, and is use-case specific. In research, it is typically set to acquire the best performance possible. Some analyze  $\theta$  as a range of values, for a more complete characterization is often obtainable with

evaluation curves that inherently show performance trade-offs. However, the held-out validation and test sets typically share data distributions as a single source partitioned into subsets (*i.e.*, train, validation, test). Regardless, the decision boundary in score space that maximizes the performance is transferred to the pin-point (*i.e.*, 1D) decision boundary - in our case, with cosine similarity, the value is a floating point value spanning [0, 1].

#### 8.4.1.2 Feature alignment

Domain  $\mathcal{D}$  can be represented by the tuple  $\mathcal{D} = \{(\mathbf{x}_i, y_i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^N$ , with  $\mathcal{X}$  and  $\mathcal{Y}$  representing the input feature space and output label space, respectively. The objective of FR algorithms is to learn a mapping function (*i.e.*, an hypothesis):  $\eta : \mathcal{X} \rightarrow \mathcal{Y}$ , which assigns each sample vector with a semantic identity label.

In domain adaptation, a model is trained on source data and deployed on target data, where an abundance of paired data is available to train a model for a task similar to that of the target. Mathematically, the labeled source domain  $\mathcal{D}_S$  and the unlabeled target domain  $\mathcal{D}_T$  can be denoted as  $\mathcal{D}_S = \{(\mathbf{x}_i^s, y_i^s) \in \mathcal{X}_S \times \mathcal{Y}_S\}_{i=1}^{N_S}$  and  $\mathcal{D}_T = \{\mathbf{x}_i^t \in \mathcal{X}_T\}_{i=1}^{N_T}$  with the sample count  $N_S = |\mathcal{D}_S|$  and  $N_T = |\mathcal{D}_T|$  corresponding to the  $i$ -th sample (*i.e.*,  $\mathbf{x}_i \in \mathbb{R}^d$ ) and label (*i.e.*,  $y_i \in \{1, \dots, K\}$ ).  $\mathcal{D}_S$  and  $\mathcal{D}_T$  are further defined by tasks  $\mathcal{T}_S$  and  $\mathcal{T}_T$ , which is indicative of the exact label type(s) and the specific  $K$  classes of interest. The goal is to learn an objective  $\eta_S : \mathcal{X}_S \rightarrow \mathcal{Y}_S$ , and then transfer to target  $\mathcal{D}_T$  for  $\mathcal{T}_T$ . By this, knowledge is leveraged from both  $\mathcal{D}_S$  for  $\mathcal{D}_T$ , with the goal of obtaining  $\eta_T$ . Since such two domains share different marginal distributions, *i.e.*,  $p(\mathbf{x}^s) \neq p(\mathbf{x}^t)$ , as well as distinct conditional distributions, *i.e.*,  $p(y^t|\mathbf{x}^s) \neq p(y^t|\mathbf{x}^t)$ , the model trained only by the labeled source domain samples usually performs poorly on the unlabeled target domain. A typical solution towards such domain gap is to learn a model  $f$  that aligns the features in a shared subspace:  $p(f(\mathbf{x}^s)) \approx p(f(\mathbf{x}^t))$ .

#### 8.4.2 Proposed framework

Given samples with identity and subgroup labels –  $\mathcal{D} = \{\mathbf{x}_i, y_i^{id}, y_i^{att}\}_{i=1}^N$ , where  $\mathbf{x} \in \mathbb{R}^d$ ,  $y^{id} \in \{1, \dots, I\}$  and  $y^{att} \in \{1, \dots, K\}$  – are used for the two objectives of the proposed framework (Figure 8.8). Hence, we aim to learn a mapping  $\mathbf{f}_{deb} = M(\mathbf{x}, \Theta_M)$  to a lower dimensional space  $\mathbf{f}_{deb} \in \mathbb{R}^{d/2}$  that preserves identity information of the target. This is dubbed the identity loss  $\mathcal{L}_{ID}$ . We also learn to do so without subgroup information, which we call the attribute (or subgroup) loss

## CHAPTER 8. BIAS IN FACE RECOGNITION

$\mathcal{L}_{ATT}$ . The total loss (*i.e.*, the final objective  $\mathcal{L} = \mathcal{L}_{ID} + \mathcal{L}_{ATT}$ ) is formed by summing the two aforementioned loss functions:

$$\mathcal{L}_{ID} = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^I \mathbf{1}_{[k=y_i^{id}]} \log(p(y = y_i^{id} | \mathbf{x}_i)), \quad (8.3)$$

$$\mathcal{L}_{ATT} = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K \mathbf{1}_{[k=y_i^{att}]} \log(p(y = y_i^{att} | \mathbf{x}_i)), \quad (8.4)$$

where  $p(y = y_i^{id} | \mathbf{x}_i)$  and  $p(y = y_i^{att} | \mathbf{x}_i)$  represent the softmax conditional probability of its identity and attribute.

We added  $\mathcal{L}_{ATT}$  to debias the features to remove variation in scores that were previously handled with a variable threshold. Further, a byproduct are these features that preserve identity information without containing knowledge of subgroup— a critical concern in the privacy and protection of biometric data.

There are three groups of parameters (*i.e.*,  $\Theta_M$ ,  $\Theta_{ID}$  and  $\Theta_{ATT}$ ) required to be optimized by the objective (Figure 8.8). Both the identity classifier  $C_{ID}$  and attribute classifier  $C_{ATT}$  are used to find a feature space that remains accurate to identity and not for subgroup by minimizing the empirical risk of  $\mathcal{L}_{ID}$  and  $\mathcal{L}_{ATT}$ :

$$\Theta_{ID}^* = \arg \min_{\Theta_{ID}} \mathcal{L}_{ID}, \quad (8.5)$$

$$\Theta_{ATT}^* = \arg \min_{\Theta_{ATT}} \mathcal{L}_{ATT}. \quad (8.6)$$

It is important to note that the task of  $\mathcal{L}_{ATT}$  is to be incorrect (*i.e.*, learn a mapping that contains no knowledge of subgroup). Thus, a gradient reversal layer [13] that acts as the identity during the forward pass, while inverting the sign of the gradient back-propagated with a scalar  $\lambda$  as the adversarial loss during training:

$$\Theta_M^* = \arg \min_{\Theta_M} -\lambda \mathcal{L}_{ATT} + \mathcal{L}_{ID}. \quad (8.7)$$

Although the proposed learning scheme is simple, it proved effective for both objectives we seek to solve. The effectiveness is clearly demonstrated in the results and analysis.

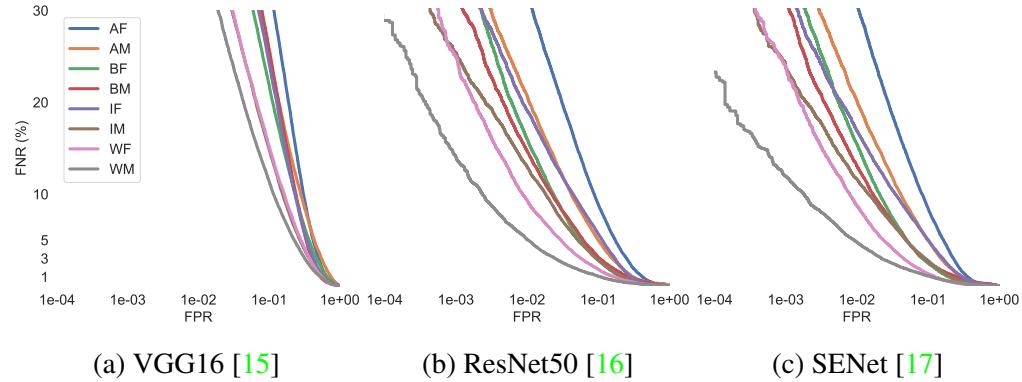


Figure 8.9: **DET curves for different CNNs.** FNR (%) (vertical) vs FPR (horizontal, log-scale) for VGG2 [14] models with different backbones (VGG16 [15], Resnet50 [16], SEnet50 [17], in that order). Lower is better. For each plot, WM is the top-performer, while AF is the worst. The ordering of the curves is roughly the same for each backbone.

## 8.5 Results and Analysis

A single CNN was used as a means to control the experiments. For this, Sphereface [48] trained on CASIA-Web [149], and evaluated on LFW [150] (%99.22 accuracy), encoded all of the faces.<sup>4</sup> As reported in [279], LFW has about 13%, 14%, 3%, and 70% ratio in Asian, Black, Indian, and White, respectively. Furthermore, CASIA-Web is even more unbalanced (again, as reported in [279]), with about 3%, 11%, 2%, and 85% for the same subgroups.

**DET analysis.** DET curves (5-fold, averaged) show per-subgroup trade-offs (Figure 8.4). Note that M performs better than F, precisely as one would expect from the tails of score-distributions for *genuine* pairs (Figure 8.3). AF and IF perform the worst.

**Score analysis.** Figure 8.3 shows score distributions for faces of the same (*i.e., genuine*) and different (*i.e., imposter*) identity, with a subgroup per SDM graph. Notice that score distributions for imposters tend to peak about zero for all subgroups, and with minimal deviation comparing modes of the different plots. On the other hand, the score distribution of the *genuine* pairs varies across subgroups in location (*i.e.*, score value) and spread (*i.e.*, overall shape). Figure 8.5 shows the confusion matrix of the subgroups. A vast majority of errors occur in the intra-subgroup. It is interesting to note that while the definition of each group based on ethnicity and race may not be

<sup>4</sup>[https://github.com/clcarwin/sphereface\\_pytorch](https://github.com/clcarwin/sphereface_pytorch)

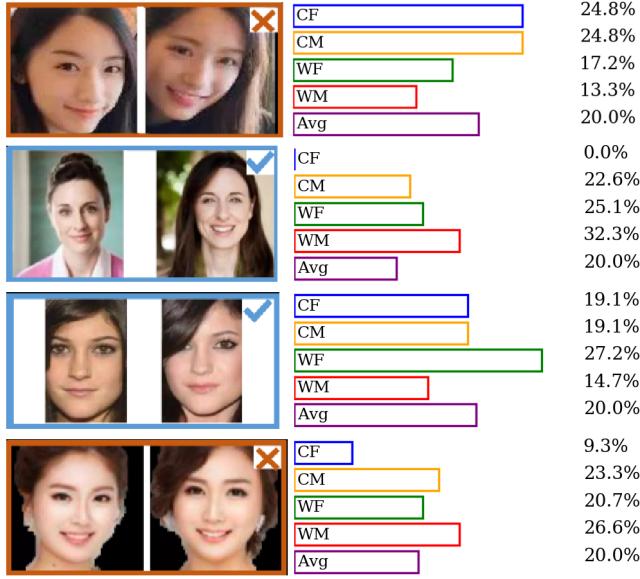


Figure 8.10: **Human assessment (qualitative).** ✓ for *match*; ✗ for *non-match*. Accuracy scores shown as bar plots. Humans are most successful at recognizing their own subgroup, with a few exceptions (*e.g.*, bottom).

crisply defined, the confusion matrix indicates that in practice, the CNN finds that the groups are effectively separate. The categories are, therefore, meaningful for FR.

The gender-based DET curves show performances with a fixed threshold (dashed line). Other curves relate similarity (lines omitted to declutter). For many FR applications, systems operate at the highest FPR allowed. The constant threshold shows that a single threshold produces different operating points (*i.e.*, FPR) across subgroups, which is undesirable. If this is the case in an industrial system, one would expect a difference in about double the FPs reported based on subgroup alone. The potential ramifications of such a bias should be considered, which it has not as of lately— even noticed in main-stream media [348, 281].

To further demonstrate the extent of the problem, we follow settings typical for systems in practice. We set the desired FPR, and then determine the percent difference, *i.e.*, desired versus actual (Figure 8.11, *top*). Furthermore, we mitigate this highly skewed phenomenon by applying subgroup-specific thresholds (*bottom*) - by this, minimal error from the desired FPR. Besides where there is a small error, the offset is balanced across subgroups.

Table 8.3: **Human assessment (quantitative)**. Subgroups listed per row (*i.e.*, human) and column (*i.e.*, image). Note, most do the best intra-subgroup (blue), and second-best (red) intra-subgroup but inter-gender. WF performs the best; WF pairs are most correctly matched.

		Image					
		(Acc, %)	CF	CM	WF	WM	Avg
Human	CF	52.9	48.0	43.8	44.7	47.4	
	CM	45.6	50.4	44.4	36.2	44.1	
	WF	44.7	43.8	57.3	48.0	48.5	
	WM	30.1	47.4	45.3	56.1	44.7	
	Avg	43.3	47.4	47.7	46.3	46.2	

**Model analysis.** Variations in optimal threshold exist across models (Figure 8.9). Like in Figure 8.4, the DET curves for three CNN-based models, each trained on VGG2 with softmax but with different backbones.<sup>5</sup> Notice similar trends across subgroups and models, which is consistent with Sphereface as well (Figure 8.4). For example, the plots generated with Sphereface and VggFace2 all have the WM curve at the bottom (*i.e.*, best) and AF on top (*i.e.*, worst). Thus, the additional CNN-based models demonstrate the same phenomena: proportional to the overall model performance, exact in the ordering of curves for each subgroup.

**Verification threshold analysis.** We seek to reduce the bias between subgroups. Such that an operating point (*i.e.*, FPR) is constant across subgroups. To accomplish that, we used a per subgroup threshold. In FV, we consider one image as the query, and all others as the test. For this, the ethnicity of the query image is assumed. We can then examine the DET curves and pick the best threshold per subgroup for a certain FPR.

We evaluated TAR for specific FAR values. As described in Section 8.3.2, the verification experiments were 5-fold, with no overlap in subject ID between folds. Results reported are averaged across folds in all cases and are shown in Table 8.5. For each subgroup, the TAR of using a global threshold is reported (upper row), as well as using the optimal per subgroup threshold (lower row).

Even for lower FAR, there are notable improvements, often of the order of 1%, which can be challenging to achieve when FAR is near  $\geq 90\%$ . More importantly, each subgroup has the desired

---

<sup>5</sup><https://github.com/rcmalli/keras-vggface>

**Table 8.4: BFW features compared to related resources.** Note, the balance across identity (ID), gender (G), and ethnicity (E). Compared with DemogPairs, BFW provides more samples per subject and subgroups per set. Also, BFW uses a single resource, VGG2. RFW; on the other hand, supports a different task (*i.e.*, subgroup classification). Furthermore, RFW and FairFace focus on race-distribution without the support of identity labels.

Database		Number of			Balanced Labels		
Name	Source	Faces	IDs	Subgroups	ID	E	G
RFW [279]	MS-Celeb-1M	≈80,000	≈12,000	4	✗	✓	✗
DemogPairs [326]	CASIA-W, VGG (+2)	10,800	600	6	✓	✓	✓
FairFace [327]	Flickr, Twitter, Web	108,000	–	10	✗	✓	✓
BFW (ours) [25]	VGG2	20,000	800	8	✓	✓	✓

FPR, so that substantial differences in FPR will remain unfounded. We experimented with ethnicity estimators on both the query and test image, which yielded similar results to those reported here.

**Human evaluation analysis.** Subjects of a subgroup likely have mostly been exposed to others of the same (Table 8.3 and Figure 8.10). Therefore, it is expected they would be best at labeling their own, similar to the same ethnicity, but another gender. Our findings concur. Each subgroup is best at labeling their type, and then second best at labeling the same ethnicity but opposite sex.

Interestingly, each group of images is best tagged by the corresponding subgroup, with the second-best having the same ethnicity and opposite gender. On average, subgroups are comparable at labeling images. Analogous to the FR system, performance ratings differ when analyzing within and across subgroups. In other words, performance on BFW improved with subgroup-specific thresholds. Similarly, humans tend to better recognize individuals by facial cues of the same or similar demographics. Put differently, as the recognition performances drop with a global threshold optimized for one subgroup and deployed on another, human performance tends to fall when across subgroups (*i.e.*, performances drop for less familiar subgroups).

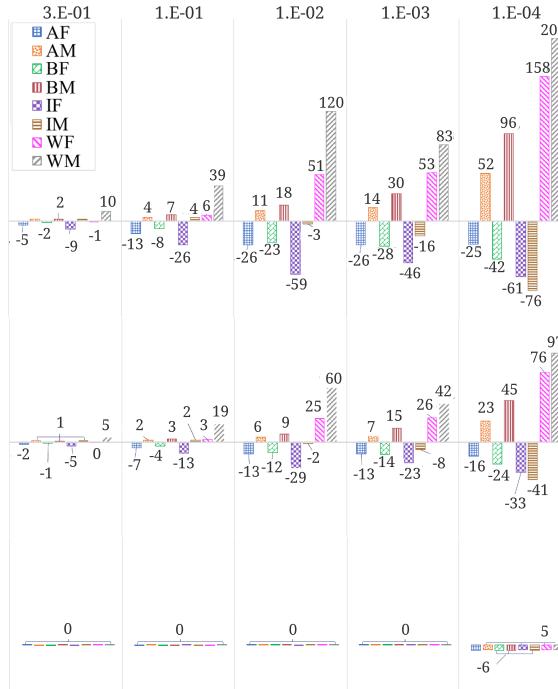


Figure 8.11: **Percent difference from intended FPR per subgroup.** *Top:* global threshold ( $t_g$ ) yields a FPR that spans up to 200% the intended (*i.e.*, WM for 1e-4); the F subgroups tend to perform worse than intended for all, while M’s overshoot the intended performances besides IM at FPR=1e-4. *Bottom:* Subgroup-specific (or optimal) thresholds  $t_o$  reduces the difference closer to zero. Furthermore, the proposed method (*middle*), which does not assume knowledge of attribute information at inference like for  $t_o$ , clearly mitigates the issue of the inconsistencies in the true versus reported FPR. Similar to the results in Table 8.5, the variations are nearly halved: the percent difference for subgroups is more balanced using the adapted features versus the baseline.

## 8.6 Experiments

Two sets of experiments are done to demonstrate the effectiveness of the proposed using our balanced BFW [25]. First, we evaluate verification performance using debiasing. Specifically, we compare the reported results compared to the actual results per subgroup. Then, for the privacy preserving claim, we compare the performance of models trained on top of debiased features  $f_{deb}$  with those of the original features  $f_{in}$ . For each, we present the problem statement, metrics and settings, and analysis. Finally, we do an ablation study to show the performance on the renowned LFW benchmark [150].

### 8.6.1 Common settings

The baseline (*i.e.*,  $f_{in}$ ) were encoded using Arcface [172] (*i.e.*, ResNet-34).<sup>6</sup> MS1M [196] was the train set, providing about 5.8M faces for 85k subjects. We prepared the faces using MTCNN [102] to detect five facial landmarks. A similarity transformation then was applied to align the face by the five detected landmarks, from which we cropped and resized each to  $96 \times 112$ . The RGB (*i.e.*, pixel values of  $[0, 255]$ ) were normalized by centering about 0 (*i.e.*, subtracting 127.5) and then standardizing (*i.e.*, dividing by 128); encodings were later L2 normalized [255]. The batch size was 200, and a stochastic gradient descent optimizer with a momentum 0.9, weight decay 5e-4, and learning rate started at 0.1 and decreased by a factor of 10 two times when the error leveled. The choice in these settings was made based on Arcface being among the best performing FR deep models to date, and as it has become a popular choice for an *off-the-shelf* option for face recognition technology and applications.

For all experiments we used our BFW dataset (Section 5.2): the debias and privacy-based experiments use the predefined five-folds<sup>7</sup>; the ablation study on LFW uses all of BFW to train M (Figure 8.8). As mentioned, BFW was built using data of VGG2, and there exist no overlap with CASIA-Webface and LFW used to train the face encoder and for the ablation, respectfully.

### 8.6.2 Debias experiment

Typical FR systems are graded by the percent error - whether to a customer, prospective staff, and so forth. In other words, specialized curves, confusion matrices, and other metrics are not always the best way to communicate system performances to non-technical audiences. It is better to discuss ratings in a manner that is more globally understood, and more comprehensible with respect to the use-case. A prime example is to share the error rate per number correctly. For instance, claiming that a system predicts 1 FP per 10,000 predictions. However, such an approximation can be hazardous, for it is inherently misleading. To demonstrate this, we ask the following questions. *Does this claim hold true for different demographics? Does this rating depend on the types of faces - does it hold for all males or all females?* Thus, if we set our system according to a desired FAR, would the claim be fair regardless of demographics (*i.e.*, subgroups) of population.

The aforementioned questions were central to our previous work [25]. We found the answer to these questions to be clear - *No, the report FAR is not true when analyzed per subgroup.*

---

<sup>6</sup><https://github.com/deepinsight/insightface>

<sup>7</sup><https://github.com/visionjo/facerec-bias-bfw>

We found when comparing the FAR values (*i.e.*, the reported-to-the-actual), the values drastically deviate from the reported average when the score threshold is fixed for all subgroups. Furthermore, demographic-specific thresholds, meaning an assumption that demographic information is known prior, proved to mitigate the problem. However, prior knowledge of demographic, although plausible (*e.g.*, identifying a known subject on a *black list*), it is a strong assumption that limits the practical uses for which it could be deployed. To extend our prior work, we propose a debiasing scheme to reduce the differences between the reported and actual. In other words, we set out to make the claim in reported false rates to be fair across all involved demographics.

### 8.6.2.1 Metrics and settings

TAR and FAR are used to examine the trade-off in the confusions that is dependent on the choice of threshold discussed earlier. Specifically, we look at actual TAR scores as a function of desired FAR. Furthermore, we compute the following metric, the percent difference of the true and reported FAR values (*i.e.*, an average score is targeted). So we ask, *how well do the different subgroups compare to that average?* Specifically,

$$\% \text{ Error}(l) = \left( \frac{\text{FAR}(l)_{\text{reported}} - \text{FAR}(l)_{\text{actual}}}{\text{FAR}(l)_{\text{reported}}} \right) * 100\% \quad (8.8)$$

### 8.6.2.2 Analysis

The proposed balances results in a way that significantly boosts the TAR at a given FAR. The percent difference between in reported to actual FAR score implies a more fair representation has been acquired Figure 8.11: left-most are the percent differences using a single threshold and the right-most is using a variable threshold (*i.e.*, results of [25]). The proposed adaptation scheme did not only preserve performances, and with a slight boost (Table 8.5), but comparing the actual to the shows a smoothing out in deviations.

Several hard positives that were incorrectly classified by the baseline but correctly identified by the proposed are shown in Figure 8.12. The set was selected for having scores closest to the global threshold of the baseline. Also, the sample pairs shown were correctly matched by the proposed. Notice the quality of one or more of the faces in each pairs is often low-resolution; additionally, extreme pose differences between faces of each pair also is common. These challenges are overcome by the proposed scheme - mitigating the issue of bias boosts results, and there displayed are several of the pairs that went from falsely being rejected to correctly being accepted.

**Table 8.5: True Acceptance Rate (FAR) for various False Acceptance Rate (FAR).** TAR scores for a global threshold (top), the proposed debiasing transformation (middle), optimal threshold (bottom). Higher is better. The standard deviation from the average is shown to demonstrate the standard error comparing the reported (*i.e.*, average) to the subgroup-specific scores. The proposed recovers most of the loss from using a global threshold rather than a per-subgroup threshold.

	FAR	0.3	0.1	0.01	0.001	0.0001
<b>AF</b>		0.990	0.867	0.516	0.470	0.465
		0.996	0.874	0.521	0.475	0.470
		1.000	0.882	0.524	0.478	0.474
<b>AM</b>		0.994	0.883	0.529	0.482	0.477
		0.996	0.886	0.531	0.484	0.479
		1.000	0.890	0.533	0.486	0.482
<b>BF</b>		0.991	0.870	0.524	0.479	0.473
		0.995	0.875	0.527	0.481	0.476
		1.000	0.879	0.530	0.484	0.480
<b>BM</b>		0.992	0.881	0.526	0.480	0.474
		0.995	0.886	0.529	0.483	0.478
		1.000	0.891	0.532	0.485	0.480
<b>IF</b>		0.996	0.881	0.532	0.486	0.481
		0.998	0.883	0.533	0.487	0.483
		1.000	0.884	0.534	0.488	0.484
<b>IM</b>		0.997	0.895	0.533	0.485	0.479
		0.998	0.897	0.534	0.486	0.480
		1.000	0.898	0.535	0.486	0.481
<b>WF</b>		0.988	0.878	0.517	0.469	0.464
		0.992	0.884	0.522	0.472	0.468
		1.000	0.894	0.526	0.478	0.474
<b>WM</b>		0.989	0.896	0.527	0.476	0.470
		0.996	0.901	0.530	0.479	0.474
		1.000	0.910	0.535	0.483	0.478
<b>Std.</b>		0.030	0.010	0.006	0.006	0.006
<b>Dev.</b>		0.002	0.009	0.005	0.005	0.005
<b>Avg.</b>		0.992	0.881	0.526	0.478	0.473
		0.998	0.886	0.528	0.481	0.476
		1.000	0.891	0.531	0.483	0.479

### 8.6.3 Privacy preserving experiment

As mentioned, our objective was two-fold. First, we aimed to preserve identity information while debiasing the facial features, as demonstrated in the prior experiment. But then, secondly, our objective function injected the reverse gradient as part of the loss to force the embedding to be unable to classify subgroups (*i.e.*, demographics). In other words, another benefit of the proposed debiasing scheme is that it rids the facial encodings of demographic information. This, in itself, is of interest in problems of privacy and protection - ideally, face encodings, which often are the only representation

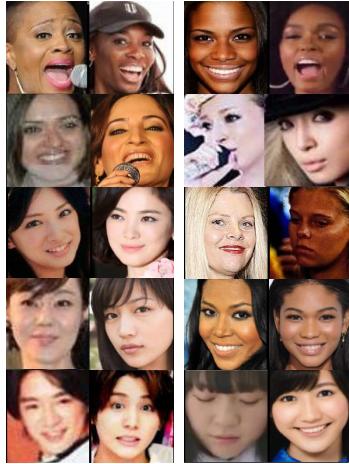


Figure 8.12: **Sample pairs of hard positives.** Pairs incorrectly classified by the baseline and correctly matched by the proposed.

of face information available at the system level, will not include attribute information, like gender or ethnicity. Hence, we aim for the inability to learn subgroup classifiers on top of the features as a means to show the demographic information has been reduced.

To show the effectiveness of the proposed in removing subgroup information, we train a MLP to classify subgroups on top of the features. We can then measure the amount of information present in the face representation [296].

The MLP designed in Keras as follows: three fully-connected layers of size 512, 512, and 256 fed into the output fully-connected layer (*i.e.*, size 8 for the 8 subgroup classes). The first three layers were separated by ReLU activation and drop-out [349] (*i.e.*, probability of 0.5) while only dropout (again, 0.5) was placed prior to the output softmax layer. A categorical crossentropy loss with Adam [261] set with a learning rate of 1e-3 was used to train.

### 8.6.3.1 Metrics and settings

We will examine the overall accuracy of the subgroup classifiers by use of a confusion matrix. Specifically, we will look at how often each subgroup was predicted correct and, when incorrect, the percentage it was mistaken for the others. The confusion was generated by averaging across the five folds.

Besides accuracy and confusions, we examine the precision and recall for each of the subgroups, along with the overall average. Precision, a measure of time correct when the prediction

Table 8.6: **Subgroup classification results.** The baseline and proposed are on the left and right columns, respectively. Note that the columns on the right have lower scores as intended.

	<b>Precision</b>		<b>Recall</b>		<b>F1</b>	
AF	0.962	0.734	0.927	0.852	0.943	0.788
AM	0.864	0.707	0.974	0.730	0.915	0.717
BF	0.940	0.655	0.924	0.644	0.932	0.647
BM	0.961	0.644	0.962	0.668	0.961	0.653
IF	0.961	0.641	0.935	0.649	0.948	0.644
IM	0.898	0.519	0.902	0.589	0.898	0.550
WF	0.934	0.554	0.970	0.547	0.951	0.549
WM	0.943	0.524	0.848	0.317	0.892	0.392
Average	0.933	0.622	0.930	0.624	0.930	0.617

assumed to be true, is calculated as follows:

$$P(l) = \frac{TP}{TP + FP}, \quad (8.9)$$

where the AP is the mean of all subgroups  $l \in L$  (*i.e.*,  $|L| = P_L$ ) for a given TPR.

The recall R, the ratio of the number of predicted-to-actual positive samples, is found as

$$R(l) = \frac{TP}{TP + FN}. \quad (8.10)$$

This complements the confusion by allowing the specificity and sensitivity of the subgroups to also be examined. Nonetheless, there are inherent trade-offs between P and R. This motivates the  $F_1$ -score [350], which fuses P and R as the harmonic mean,

$$F_1 = 2 * \frac{P * R}{P + R}. \quad (8.11)$$

### 8.6.3.2 Analysis

We demonstrated that identity knowledge is preserved (Table 8.5), and now we show the other benefits in privacy. The results clearly show that the privacy preserving claim is accurate, leading to a 30% drop in the ability to predict gender and ethnicity from the encodings (Table 8.6).

**CHAPTER 8. BIAS IN FACE RECOGNITION**

	<b>AF</b>	7.0	0.0	0.0	0.1	0.0	0.2	0.0
<b>AM</b>	1.6	<b>97.4</b>	0.1	0.1	0.0	0.2	0.0	0.7
<b>BF</b>	0.8	0.0	<b>92.4</b>	2.8	0.9	0.0	3.0	0.0
<b>BM</b>	0.0	0.0	2.0	<b>96.2</b>	0.0	1.6	0.0	0.2
<b>IF</b>	0.9	0.0	3.0	0.0	<b>93.5</b>	0.4	2.2	0.0
<b>IM</b>	0.0	3.6	0.0	0.8	2.0	<b>90.2</b>	0.0	3.3
<b>WF</b>	0.4	0.4	0.8	0.0	0.4	0.0	<b>97.0</b>	1.1
<b>WM</b>	0.0	4.6	0.0	0.1	0.3	8.6	1.6	<b>84.8</b>
	<b>AF</b>	<b>AM</b>	<b>BF</b>	<b>BM</b>	<b>IF</b>	<b>IM</b>	<b>WF</b>	<b>WM</b>

(a) Baseline subgroup classifier.

	<b>AF</b>	8.8	1.1	0.3	1.7	0.6	1.5	1.0
<b>AM</b>	10.2	<b>73.0</b>	3.1	1.2	1.6	3.8	4.6	2.6
<b>BF</b>	3.5	3.8	<b>64.4</b>	11.4	5.4	3.2	6.5	1.8
<b>BM</b>	1.5	2.1	13.3	<b>66.8</b>	4.5	5.7	4.0	2.1
<b>IF</b>	2.2	2.3	4.6	5.7	<b>64.9</b>	13.9	4.2	2.2
<b>IM</b>	1.7	3.7	1.4	5.4	9.9	<b>58.9</b>	8.4	10.6
<b>WF</b>	1.4	4.5	5.2	7.9	5.8	11.7	<b>54.7</b>	8.8
<b>WM</b>	10.4	5.4	6.1	6.0	7.9	17.1	15.4	<b>31.7</b>
	<b>AF</b>	<b>AM</b>	<b>BF</b>	<b>BM</b>	<b>IF</b>	<b>IM</b>	<b>WF</b>	<b>WM</b>

(b) Our subgroup classifier.

Figure 8.13: **Subgroup confusion matrix.** Comparison of accuracy in classifying and misclassifying the subgroups. Notice the (b) performs significantly worse than (a) as intended.

Hence, predictive power of all subgroups dropped significantly. Furthermore, the drop in performance is sufficient enough to make the claim that the predictions are now unreliable. Honing in on the specifics, it is interesting to note that the subgroups for which the baseline were most in favor of are hindered the most from the debias scheme. In other words, WM and WF drop the most, while the AM and AF drop the least. All the while the same trends in confusions propagate from the baseline to the proposed results. For instance, WM are mostly confused as IM originally, and then again in the case for the debiased features. The same holds for the opposite sex in all cases.

Next, we examine the confusions for the different subgroups before and after debiasing the face features (Figure 8.13). As established, the baseline contains more subgroup knowledge—a model can learn on top. When trained and evaluated on BFW, the baseline performs better on F subgroups. This differs from the norm where M are a majority of the data. To the contrary, the WM are inferior in performances to all subgroups in either case.

#### 8.6.4 Ablation study

To check the effectiveness of the proposed scheme we train M using the entire BFW dataset and deploy on the well-known LFW benchmark. We note that the training dataset that we employ is significantly smaller than that used by SOTA networks trained to achieve high performance on LFW, which employ the MS1MV2 dataset, which contains 5.8M images of 85k identities. By contrast, even though we initialize our network starting with features learnt on MS1MV2, we train on a small dataset of 20k images of 800 subjects, which is two orders of magnitude smaller. Although we train The current SOTA with 99.8% verification accuracy, while the proposed scheme reaches its best score of 95.2% after 5 epochs before dropping off and then leveling out around 81% (Figure 8.14). Clearly, the benefits of privacy and debiasing are hindered on unbalanced data (*i.e.*, LFW is made up of about 85% WM). Furthermore, we optimized M by choosing the epoch with the best performance prior to the drop off. Future steps could be improving the proposed when transferring to unbalanced sets with a way to detect the optimal training settings.

### 8.7 Discussion

We introduced the Balanced Faces In the Wild (BFW) dataset with eight subgroups (*i.e.*, different gender and ethnicity) for which data is split evenly across. With this, we provide evidence that the subgroups we chose and formed is meaningful, *i.e.*, the FR algorithm rarely makes mistakes

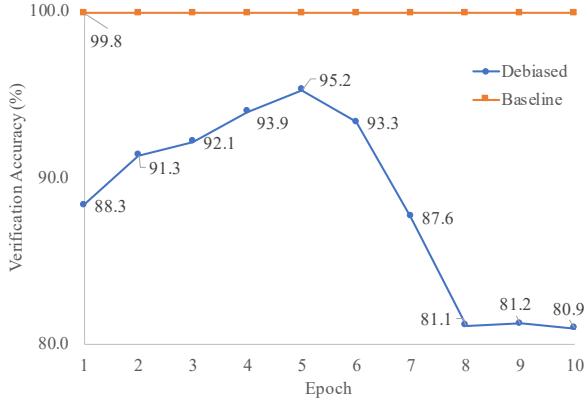


Figure 8.14: **Accuracy on LFW benchmark.** The proposed approaches the performance of the baseline before dropping off.

across subgroups. We used *off-the-shelf* CNNs, hypothesizing that these SOTA CNN suffers from bias because of the imbalanced train-set. Once established that the results do suffer from problems of bias, we observed that the same threshold across ethnic and gender subgroups leads to differences in the FPR up to a factor of two. Also, we clearly showed notable percent differences in ratings across subgroups. Furthermore, we ameliorate these differences with a per-subgroup threshold, leveling out FPR, and achieving a higher TPR. We hypothesized that most humans grew among more than their own demographic and, therefore, effectively learn from imbalanced datasets. In essence, a human evaluation validated that humans are biased, as most recognized their personal demographic best. This research, along with the data and resources, are extendable in vast ways. Thus, this is just a sliver of the larger problem of bias in ML.

We show a bias for subgroups in FR caused by the selection of a single threshold. Previously, a subgroup-specific threshold was proposed as a solution for the case in which such knowledge (*i.e.*, subgroup information) is accessible prior. Inspired by works in feature alignment and domain adaptation we propose a scheme to mitigate the bias problem. We learn a lower dimensional mapping that preserves identity and removes knowledge of subgroup. The encodings balance the performance across subgroups, while boosting the overall accuracy. Using a single threshold, the difference between actual and average performance across subgroups is reduced. Furthermore, the resulting encodings hold reduced knowledge of subgroups, increasing privacy.

# Chapter 9

## Discussion

We now wrap-up the dissertation with a discussion on the broader impacts, future directions in practice and in research, and, finally, a concluding section. To keep things focused on the main topic, the closing discussion is on topics of automatic kinship understanding and modeling. Hence, the other subtopics were handled in a more open and close manner (*i.e.*, provided as preliminary or secondary knowledge in nearly standalone chapters), we chose to omit from this final chapter as to appropriately conclude the majority of the research tailored to the kin-based topics.

### 9.1 Broader impacts

We believe that, collectively, our greatest contribution for automatic kinship recognition was the labeled data, *i.e.*, the FIW dataset.<sup>1</sup> Also, the task evaluations released with clearly defined problem statements, task protocols, and benchmark reports provide additional structure for researchers to follow as we continue to push the envelope of capabilities and technology developed as solutions for the problems. Thus, by laying the groundwork for others to get started and incorporate modern-day, data driven modeling ideas, we hope the trend continues in that attention for this problem continues to climb as our automatic face analysis and modeling capabilities progress as well.

The fourth RFIW gained fair attention—T1 (verification) saw the most; T2 (tri-subject) and T3 (search and retrieval), both supported for the first time, are more complex and practically motivated than the classic T1. The broader impact spans greater than current tasks in application (*e.g.*, generative-based tasks [10, 219]) and experimental settings (*e.g.*, with privacy a concern [241]). RFIW met the difficulty and practicality of today; the question how best to formulate the problem

---

<sup>1</sup>More information, downloads, and publications are on the project page, <https://web.northeastern.edu/smilelab/fiw/>.

## CHAPTER 9. DISCUSSION

remains an open research question. As such, this survey aims to provide a stronghold on the laboratory-style evaluations as seen appropriate in the modern day.

### 9.2 Future work

It is an exciting and opportune time for kin-based problems for researchers and practitioners alike. For starters, there is a lot of room for improving SOTA, and even the experiments (*i.e.*, design, purpose, and extent). For instance, incorporating additional label types (*i.e.*, other soft attributes like expression, age, and ethnicity), different data splits and protocols (*e.g.*, given a father, daughter, and grandparents from the side of the mother, determine the mother), and practical use-cases (*e.g.*, automate family photo-album creation). Generative-based tasks also hold promise in directions to take next: whether improved predictive capability of a child’s face - provided a pair of parents, or a more fine-grained view of predicting any node in a family tree - provided samples of all other family members - then the room for improvement and potential for growth is furthered.

**Fairness.** A few recent attempts have been made by researchers to address fair AI and transparency in kinship understanding. For instance, the latest version of RFIW, supplemental to this survey, FIW is now supported with a *datasheet*: “datasheets for datasets” [351].<sup>2</sup> The motivation of datasheets is to promote transparency and, thus, to minimize the doubt from unknown biases that come and are inherited by publicly available data resources. Specifically, *datasheets* completely spec-out the data (*e.g.*, motivation, composition, collection process, preprocessing, updates, legal and ethical considerations). There are other methods for transparency that have been recently proposed with a similar motivation as “datasheets for datasets”, such as *fact sheets* [352] and *model cards* [353]. Nonetheless, we found that the format and motivation of “datasheets for datasets” as the best for FIW. So, this was used to record and archive data specifications.

**Privacy.** As is the case for many ML tasks, privacy has motivated researchers. Recently, Kumar *et al.* proposed using a GNN to first achieve SOTA in family classification, and to then add noise to encrypt the data, and demonstrating that a variant of the model safely encapsulates the learned knowledge (*i.e.*, an ability to accurately deceiver) [241].

---

<sup>2</sup>The datasheet for FIW is available online, [https://web.northeastern.edu/smilelab/fiw/fiw\\_ds.pdf](https://web.northeastern.edu/smilelab/fiw/fiw_ds.pdf).

## CHAPTER 9. DISCUSSION

**Social and cultural.** A near radical piece of its time, Goode [354] surveyed family structure as more of a complex system than the ‘conjugal family form’ of many traditional cultures (*e.g.*, Western, Chinese, Arab). Hence, we currently look at sets of persons as being either related or not related—this does not account for the realistic setting that would be faced in the modern world. For instance, tri-subject pairs are structured by using a parent pair and a child (*i.e.*, using evidence of both parents)—a scenario that is certainly a step in the right direction, as parents are often inferred through records on marriage and offspring. However, families are dynamic in many modern cultures—step siblings and parents are common. Considering the setting of tri-subject: what happens if the father is true, but the mother is not, or vice versa— a concept that propagates to all levels of the problem, and especially when considering complete family trees with connections to in-laws. Thus, the problem remains: how to best weight (*i.e.*, fuse [355]) different relationship types. Even simple questions have soft, varying solutions [356] like *Do we look more like our father?*

**Feature fusion.** Still today, the underlying question remains. *How to best fuse prior knowledge?* For instance, in tri-subject verification, the fusion of the features from the two parents. Flipping this very problem around (*i.e.*, given parents, generate the child’s face), the question of feature fusion is still prominent. Looking ahead at attempts to solve the fine-grained problem of populating family trees, regardless if viewed as discriminative or generative, the question remains: *how to best leverage prior knowledge of additional family members relatively of different types and degrees?*

Although the number of methods is great—whether metric-learning, deep features, a variant of both—most recent attempts only differ in the broad sense. Bottom-line, successes in all tasks have been tributes of systems based on a Siamese network(s) that encodes inputs from image-to-feature space. The feature space learned typically differs in the point and method of fusion. Specifically, paired samples are usually split evenly (*i.e.*, the number of pair-types of type *KIN* and *NON-KIN* for each relationship type is split *fifty-fifty*). Provided a Siamese network, often pre-trained on auxiliary face recognition dataset, act as face encoders. In order to transform from feature-to-score space, either a metric, fusion technique, or both are applied—this tends to be where methods differ, yet the same conventional coarse system holds (*i.e.*, Figure 4.5). In summary, it is the Siamese net to encode faces, followed by some means of feature-fusion that are stove-piped to a metric or learning objective. Hence, some relevant aspects of such a system produce current SOTA from which we had drawn conclusions, and especially in identifying research trends and open issues. We consider the most relevant among aspects for achieving effective systems as follows: (1) effective method for fusion; (2) representation that considers the relationship’s direction; (3) detecting other attributes (*e.g.*, age

## CHAPTER 9. DISCUSSION

and gender) and knowledge of the higher-level scene (*e.g.*, face detected in picture with car styles that hint the picture was taken in the 1970s).

**Multimodal data.** Let us consider other signals that can define visual data; let us consider other label types for faces that could also enhance performance. For instance, expressions and mannerisms are often similar for parent and child (*e.g.*, *they have the mother's smile*). More complex dynamics for individual expressions and mannerisms can be effectively captured in video data [357]. Hence, added knowledge that complements the visual information has proven useful for boosting kin-based recognition ratings (*e.g.*, 3D facial images[107], voice [137], MM [19]).

**Family synthesis.** The existing technology for synthesizing family members is still immature and generalization remains unsolved. A system should accept one-to-many members of a family tree and synthesize the appearances of the desired relative type. Still, many questions remain: how to fuse, handle dynamic inputs, the optimal way to reason about a family tree, and more.

### 9.3 Conclusion

Upon completing the dissertation, and looking back at years of research, we see decent progress in automatic kinship recognition, along with an increase in attention, an improvement in the methodology, advancement in available data and resources, and improved benchmarks (*i.e.*, both in protocol and proposed SOTA). This dissertation contributed in each of the aforementioned ways: we improved models for kin-specific face-based recognition tasks; we attracted researchers by organizing several workshops, challenges, tutorials, etc.; we released our large-scale FIW data collection, along with the data splits and benchmark scores for public use. Besides discriminative tasks (*i.e.*, recognition), for which we supported various tasks to span different use-cases, we also explored generative-based problems. Specifically, we were the first to attempt the appearance of off-spring from a pair of parents as the input (*i.e.*, not a single parent, but with two-to-one mapping so more information was available for inference). Also, we teased the use of multimedia in kinship recognition with our FIW-MM—the benefits of the added modalities were clearly demonstrated. Finally, we recently delivered a comprehensive survey on the topic as means of establishing a stronghold on the state of technology after a decade of research— a single resource to reference for experimental protocols, benchmark ratings, and a high-level view of the latest-and-greatest methods. We believe the research for this dissertation was well spent for the advancing of automatic kinship understanding.

# Bibliography

- [1] H. Wang, Y. Wang, Z. Zhou, X. Ji, D. Gong, J. Zhou, Z. Li, and W. Liu, “Cosface: Large margin cosine loss for deep face recognition,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [2] A.-W. Harzing, *The publish or perish book*. Tarma Software Research Pty Limited, 2010.
- [3] J. Ramos *et al.*, “Using tf-idf to determine word relevance in document queries,” in *International Conference on Machine Learning*, 2003.
- [4] J. P. Robinson, M. Shao, Y. Wu, and Y. Fu, “Families in the wild (fiw): Large-scale kinship image database and benchmarks,” in *ACM Conference on Multimedia (ACMMM)*, 2016.
- [5] J. P. Robinson, M. Shao, H. Zhao, Y. Wu, T. Gillis, and Y. Fu, “Recognizing families in the wild (rfiw),” in *RFIW at ACM MM*, 2017.
- [6] J. P. Robinson, M. Shao, Y. Wu, H. Liu, T. Gillis, and Y. Fu, “Visual kinship recognition of families in the wild,” *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, 2018.
- [7] J. P. Robinson, Y. Yin, Z. Khan, M. Shao, S. Xia, M. Stopa, S. Timoner, M. A. Turk, R. Chellappa, and Y. Fu, “Recognizing families in the wild (rfiw): The 4th edition,” *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2020.
- [8] A. Dehghan, E. Ortiz, R. Villegas, and M. Shah, “Who do i look like? determining parent-offspring resemblance via gated autoencoders,” in *Computer Vision and Pattern Recognition*, 2014.
- [9] R. Fang, A. Gallagher, T. Chen, and A. Loui, “Kinship classification by modeling facial feature heredity,” in *International Conference on Image Processing (ICIP)*. IEEE, 2013.

## BIBLIOGRAPHY

- [10] P. Gao, S. Xia, J. P. Robinson, J. Zhang, C. Xia, M. Shao, and Y. Fu, “What will your child look like? dna-net: Age and gender aware kin face synthesizer,” *CoRR arXiv:1911.07014*, 2019.
- [11] S. Wang, Z. Ding, and Y. Fu, “Cross-generation kinship verification with sparse discriminative metric,” *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 41, no. 11, 2018.
- [12] M. Dawson, A. Zisserman, and C. Nellåker, “From same photo: Cheating on visual kinship challenges,” in *Asian Conference on Computer Vision*. Springer, 2018.
- [13] Y. Ganin and V. Lempitsky, “Unsupervised domain adaptation by backpropagation,” in *International Conference on Machine Learning (ICML)*. PMLR, 2015, pp. 1180–1189.
- [14] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman, “Vggface2: A dataset for recognising faces across pose and age,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2018.
- [15] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *International Conference on Learning Representations (ICLR)*, 2015.
- [16] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [17] J. Hu, L. Shen, and G. Sun, “Squeeze-and-excitation networks,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [18] J. P. Robinson, M. Shao, and Y. Fu, “Visual kinship recognition: A decade in the making,” *CoRR arXiv: 2006.16033*, 2020.
- [19] J. P. Robinson, Z. Khan, Y. Yin, M. Shao, and Y. Fu, “Families in wild multimedia (fiw-mm): A multi-modal database for recognizing kinship,” *CoRR arXiv:2007.14509*, 2020.
- [20] Y. Yin, J. P. Robinson, and Y. Fu, “Multimodal in-bed pose and shape estimation under the blankets,” *CoRR arXiv:2012.06735*, 2020.
- [21] Y. Yin, J. P. Robinson, S. Jiang, Y. Bai, C. Qin, and Y. Fu, “Superfront: From low-resolution to high-resolution frontal face synthesis,” *CoRR arXiv:2012.04111*, 2020.

## BIBLIOGRAPHY

- [22] C. Zheng, S. Xia, J. P. Robinson, C. Lu, W. Wu, C. Qian, and M. Shao, “Localin reshuffle net: Toward naturally and efficiently facial image blending,” in *Asian Conference on Computer Vision*, November 2020.
- [23] Y. Yin, S. Jiang, J. P. Robinson, and Y. Fu, “Dual-attention gan for large-pose face frontalization,” *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2020.
- [24] L. Wang, B. Sun, J. P. Robinson, T. Jing, and Y. Fu, “Ev-action: Electromyography-vision multi-modal action dataset,” *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2020.
- [25] J. P. Robinson, G. Livitz, Y. Henon, C. Qin, Y. Fu, and S. Timoner, “Face recognition: too bias, or not too bias?” in *Conference on Computer Vision and Pattern Recognition Workshop*, 2020, pp. 0–1.
- [26] Y. Yin, J. Robinson, Y. Zhang, and Y. Fu, “Joint super-resolution and alignment of tiny faces,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 12 693–12 700.
- [27] W. Zhuang, Y. Wang, J. P. Robinson, C. Wang, M. Shao, Y. Fu, and S. Xia, “Towards 3d dance motion synthesis and control,” *CoRR arXiv:2006.05743*, 2020.
- [28] J. P. Robinson, Y. Li, N. Zhang, Y. Fu, and S. Tulyakov, “Laplace landmark localization,” in *IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 10 103–10 112.
- [29] Y. Wu, Z. Ding, H. Liu, J. P. Robinson, and Y. Fu, “Kinship classification through latent adaptive subspace,” in *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*. IEEE, 2018, pp. 143–149.
- [30] J. P. Robinson, M. Shao, and Y. Fu, “To recognize families in the wild: A machine vision tutorial,” in *ACM Conference on Multimedia (ACMMM)*, 2018.
- [31] S. Wang, J. P. Robinson, and Y. Fu, “Kinship verification on families in the wild with marginalized denoising metric learning,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2017.

## BIBLIOGRAPHY

- [32] J. P. Robinson and Y. Fu, “Pre-trained d-cnn models for detecting complex events in unconstrained videos,” in *SPIE Commercial+ Scientific Sensing and Imaging*. International Society for Optics and Photonics, 2016.
- [33] J. P. Robinson, E. Scott, and Y. Fu, “Neu mitll@ trecvid 2015: Multimedia event detection by deep feature learning,” in *Proceedings of TRECVID*, 2015.
- [34] R. Fang, K. D. Tang, N. Snavely, and T. Chen, “Towards computational models of kinship verification,” in *International Conference on Image Processing (ICIP)*. IEEE, 2010.
- [35] E. Learned-Miller. (2018) Elements of modern face recognition. [Online]. Available: [https://people.cs.umass.edu/~elm/Teaching/ppt/370/370\\_Face\\_Intro.pptx.pdf](https://people.cs.umass.edu/~elm/Teaching/ppt/370/370_Face_Intro.pptx.pdf)
- [36] P. N. Belhumeur, J. P. Hespanha, and D. J. Kriegman, “Eigenfaces vs. fisherfaces: Recognition using class specific linear projection,” *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 19, no. 7, pp. 711–720, 1997.
- [37] A. S. Georghiades, P. N. Belhumeur, and D. J. Kriegman, “From few to many: Illumination cone models for face recognition under variable lighting and pose,” *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 23, no. 6, pp. 643–660, 2001.
- [38] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” *International Conference on Learning Representations (ICLR)*, 2015.
- [39] J. Ba, V. Mnih, and K. Kavukcuoglu, “Multiple object recognition with visual attention,” *International Conference on Learning Representations (ICLR)*, 2015.
- [40] J. K. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, “Attention-based models for speech recognition,” in *Advances in Neural Information Processing Systems (NIPS)*, 2015, pp. 577–585.
- [41] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, “Graph attention networks,” *International Conference on Learning Representations (ICLR)*, 2018.
- [42] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, “Backpropagation applied to handwritten zip code recognition,” *Neural computation*, vol. 1, no. 4, pp. 541–551, 1989.

## BIBLIOGRAPHY

- [43] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems (NIPS)*, 2012, pp. 1097–1105.
- [44] J. Yang, Y.-G. Jiang, A. G. Hauptmann, and C.-W. Ngo, “Evaluating bag-of-visual-words representations in scene classification,” in *Proceedings of the international workshop on Workshop on multimedia information retrieval*, 2007, pp. 197–206.
- [45] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf, “Deepface: Closing the gap to human-level performance in face verification,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [46] O. M. Parkhi, A. Vedaldi, and A. Zisserman, “Deep face recognition,” in *British Machine Vision Conference (BMVC)*, 2015.
- [47] Y. Wen, K. Zhang, Z. Li, and Y. Qiao, “A discriminative feature learning approach for deep face recognition,” in *European Conference on Computer Vision (ECCV)*, 2016, pp. 499–515.
- [48] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song, “Spherenet: Deep hypersphere embedding for face recognition,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [49] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, “Arcface: Additive angular margin loss for deep face recognition,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [50] Y. Duan, J. Lu, and J. Zhou, “Uniformface: Learning deep equidistributed representation for face recognition,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 3415–3424.
- [51] Z. Wang, G. Wang, B. Huang, Z. Xiong, Q. Hong, H. Wu, P. Yi, K. Jiang, N. Wang, Y. Pei *et al.*, “Masked face recognition dataset and application,” *CoRR arXiv:2003.09093*, 2020.
- [52] G. Guo and N. Zhang, “A survey on deep learning based face recognition,” *Computer Vision and Image Understanding*, 2019.
- [53] I. Masi, Y. Wu, T. Hassner, and P. Natarajan, “Deep face recognition: A survey,” in *2018 31st SIBGRAPI conference on graphics, patterns and images (SIBGRAPI)*. IEEE, 2018, pp. 471–478.

## BIBLIOGRAPHY

- [54] G. Koch, R. Zemel, and R. Salakhutdinov, “Siamese neural networks for one-shot image recognition,” in *ICML deep learning workshop*, vol. 2. Lille, 2015.
- [55] J. S. Supancic, G. Rogez, Y. Yang, J. Shotton, and D. Ramanan, “Depth-based hand pose estimation: data, methods, and challenges,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 1868–1876.
- [56] T. F. Cootes and C. J. Taylor, “Active shape models—‘smart snakes’,” in *British Machine Vision Conference (BMVC)*, 1992.
- [57] T. F. Cootes, G. J. Edwards, and C. J. Taylor, “Active appearance models,” *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 23, no. 6, pp. 681–685, 2001.
- [58] R. Chellappa, C. L. Wilson, and S. Sirohey, “Human and machine recognition of faces: A survey,” *Proceedings of the IEEE*, vol. 83, no. 5, pp. 705–741, 1995.
- [59] A. Newell, K. Yang, and J. Deng, “Stacked hourglass networks for human pose estimation,” in *European Conference on Computer Vision (ECCV)*. Springer, 2016, pp. 483–499.
- [60] X. Peng, R. S. Feris, X. Wang, and D. N. Metaxas, “Red-net: A recurrent encoder–decoder network for video-based face alignment,” *International Journal of Computer Vision (IJCV)*, 2018.
- [61] J. Yang, Q. Liu, and K. Zhang, “Stacked hourglass network for robust facial landmark localisation,” in *Conference on Computer Vision and Pattern Recognition Workshop*, 2017, pp. 79–87.
- [62] S. Honari, P. Molchanov, S. Tyree, P. Vincent, C. Pal, and J. Kautz, “Improving landmark localization with semi-supervised learning,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [63] S. Ren, X. Cao, Y. Wei, and J. Sun, “Face alignment at 3000 fps via regressing local binary features,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014, pp. 1685–1692.
- [64] V. Kazemi and J. Sullivan, “One millisecond face alignment with an ensemble of regression trees,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.

## BIBLIOGRAPHY

- [65] G. Trigeorgis, P. Snape, M. A. Nicolaou, E. Antonakos, and S. Zafeiriou, “Mnemonic descent method: A recurrent process applied for end-to-end face alignment,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 4177–4187.
- [66] W. Wang, S. Tulyakov, and N. Sebe, “Recurrent convolutional face alignment,” in *Asian Conference on Computer Vision*. Springer, 2016, pp. 104–120.
- [67] ———, “Recurrent convolutional shape regression,” *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, 2018.
- [68] R. A. Güler, G. Trigeorgis, E. Antonakos, P. Snape, S. Zafeiriou, and I. Kokkinos, “Densereg: Fully convolutional dense shape regression in-the-wild,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [69] L. A. Jeni, J. F. Cohn, and T. Kanade, “Dense 3d face alignment from 2d videos in real-time,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2015.
- [70] S. Tulyakov, L. A. Jeni, J. F. Cohn, and N. Sebe, “Viewpoint-consistent 3d face alignment,” *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, 2018.
- [71] X. Zhu, Z. Lei, X. Liu, H. Shi, and S. Z. Li, “Face alignment across large poses: A 3d solution,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 146–155.
- [72] X. Dong, S.-I. Yu, X. Weng, S.-E. Wei, Y. Yang, and Y. Sheikh, “Supervision-by-registration: An unsupervised approach to improve the precision of facial landmark detectors,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 360–368.
- [73] R. Gross, I. Matthews, J. Cohn, T. Kanade, and S. Baker, “Multi-pie,” *Image and Vision Computing*, vol. 28, no. 5, 2010.
- [74] X. Zhu and D. Ramanan, “Face detection, pose estimation, and landmark localization in the wild,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2012, pp. 2879–2886.
- [75] M. Koestinger, P. Wohlhart, P. M. Roth, and H. Bischof, “Annotated facial landmarks in the wild: A large-scale, real-world database for facial landmark localization,” in *IEEE International Conference on Computer Vision (ICCV) Workshop*, 2011, pp. 2144–2151.

## BIBLIOGRAPHY

- [76] C. Sagonas, G. Tzimiropoulos, S. Zafeiriou, and M. Pantic, “300 faces in-the-wild challenge: The first facial landmark localization challenge,” in *IEEE International Conference on Computer Vision (ICCV) Workshop*, 2013.
- [77] A. Bulat and G. Tzimiropoulos, “How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks),” in *IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [78] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *Advances in Neural Information Processing Systems (NIPS)*, 2014, pp. 2672–2680.
- [79] A. Radford, L. Metz, and S. Chintala, “Unsupervised representation learning with deep convolutional generative adversarial networks,” *CoRR arXiv:1511.06434*, 2015.
- [80] M. Saito, E. Matsumoto, and S. Saito, “Temporal generative adversarial nets with singular value clipping,” in *IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [81] S. Tulyakov, M.-Y. Liu, X. Yang, and J. Kautz, “Mocogan: Decomposing motion and content for video generation,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 1526–1535.
- [82] Z. Geng, C. Cao, and S. Tulyakov, “3d guided fine-grained face manipulation,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [83] H. Zhang, T. Xu, H. Li, S. Zhang, X. Wang, X. Huang, and D. Metaxas, “Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks,” in *IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [84] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, “Unpaired image-to-image translation using cycle-consistent adversarial networks,” in *IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [85] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, G. Liu, A. Tao, J. Kautz, and B. Catanzaro, “Video-to-video synthesis,” *CoRR arXiv:1808.06601*, 2018.
- [86] A. Siarohin, S. Lathuilière, S. Tulyakov, E. Ricci, and N. Sebe, “Animating arbitrary objects via deep motion transfer,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

## BIBLIOGRAPHY

- [87] Z. Ding, Y. Guo, L. Zhang, and Y. Fu, “One-shot face recognition via generative learning,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*. IEEE, 2018, pp. 1–7.
- [88] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. Efros, and T. Darrell, “Cycada: Cycle-consistent adversarial domain adaptation,” in *International Conference on Machine Learning (ICML)*. PMLR, 2018, pp. 1989–1998.
- [89] A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, and R. Webb, “Learning from simulated and unsupervised images through adversarial training.” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [90] W. Yang, W. Ouyang, X. Wang, J. Ren, H. Li, and X. Wang, “3d human pose estimation in the wild by adversarial learning,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [91] O. Chapelle and M. Wu, “Gradient descent optimization of smoothed information retrieval metrics,” *Information retrieval*, vol. 13, no. 3, pp. 216–235, 2010.
- [92] M. D. Hoffman, D. M. Blei, C. Wang, and J. Paisley, “Stochastic variational inference,” *The Journal of Machine Learning Research*, vol. 14, no. 1, pp. 1303–1347, 2013.
- [93] P. Domingos, “A unified bias-variance decomposition,” in *International Conference on Machine Learning (ICML)*, 2000, pp. 231–238.
- [94] S. Honari, J. Yosinski, P. Vincent, and C. Pal, “Recombinator networks: Learning coarse-to-fine feature aggregation,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 5743–5752.
- [95] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, “Image-to-image translation with conditional adversarial networks,” *CoRR arXiv:1611.07004*, 2017.
- [96] X. Xiong and F. De la Torre, “Supervised descent method and its applications to face alignment,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013, pp. 532–539.
- [97] Z. Zhang, P. Luo, C. C. Loy, and X. Tang, “Facial landmark detection by deep multi-task learning,” in *European Conference on Computer Vision (ECCV)*. Springer, 2014, pp. 94–108.

## BIBLIOGRAPHY

- [98] S. Zhu, C. Li, C. Change Loy, and X. Tang, “Face alignment by coarse-to-fine shape searching,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 4998–5006.
- [99] J.-J. Lv, X. Shao, J. Xing, C. Cheng, X. Zhou *et al.*, “A deep regression architecture with two-stage re-initialization for high performance facial landmark detection.” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [100] D. Cristinacce and T. F. Cootes, “Feature detection and tracking with constrained local models.” in *British Machine Vision Conference (BMVC)*, vol. 1, no. 2. Citeseer, 2006, p. 3.
- [101] A. Nech and I. Kemelmacher-Shlizerman, “Level playing field for million scale face recognition,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [102] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao, “Joint face detection and alignment using multitask cascaded convolutional networks,” *Signal Processing Letters*, vol. 23, no. 10, pp. 1499–1503, 2016.
- [103] F. F. Furstenberg, “Kinship reconsidered: Research on a neglected topic,” *Journal of Marriage and Family*, vol. 82, no. 1, 2020.
- [104] N. Kohli, D. Yadav, M. Vatsa, R. Singh, and A. Noore, “Supervised mixed norm autoencoder for kinship verification in unconstrained videos,” *IEEE Trans. on Image Processing (TIP)*, 2018.
- [105] M. Shao, S. Xia, and Y. Fu, “Genealogical face recognition based on ub kinface database,” in *Conference on Computer Vision and Pattern Recognition Workshop*. IEEE, 2011.
- [106] S. Xia, M. Shao, and Y. Fu, “Kinship verification through transfer learning,” in *International Joint Conferences on AI (IJCAI)*, 2011.
- [107] V. Vijayan, K. W. Bowyer, P. J. Flynn, D. Huang, L. Chen, M. Hansen, O. Ocogueda, S. K. Shah, and I. A. Kakadiaris, “Twins 3d face recognition challenge,” in *2011 International Joint Conference on Biometrics (IJCB)*. IEEE, 2011, pp. 1–7.
- [108] A. G. Bottino, M. De Simone, A. Laurentini, and T. Vieira, “A new problem in face image analysis: finding kinship clues for siblings pairs,” in *Conference on Pattern Recognition Application & Methods*, 2012.

## BIBLIOGRAPHY

- [109] Y. Guo, H. Dibeklioglu, and L. van der Maaten, “Graph-based kinship recognition,” in *International Conference on Pattern Recognition (ICPR)*, 2014.
- [110] H. Dibeklioglu, A. Ali Salah, and T. Gevers, “Like father, like son: Facial expression dynamics for kinship verification,” in *IEEE International Conference on Computer Vision (ICCV)*, 2013.
- [111] E. Boutellaa, M. B. López, S. Ait-Aoudia, X. Feng, and A. Hadid, “Kinship verification from videos using spatio-temporal texture features and deep learning,” *CoRR arXiv:1708.04069*, 2017.
- [112] J. Lu, J. Hu, V. E. Liong, X. Zhou, A. Bottino, I. U. Islam, T. F. Vieira, X. Qin, X. Tan, S. Chen *et al.*, “The fg 2015 kinship verification in the wild evaluation,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, vol. 1. IEEE, 2015, pp. 1–7.
- [113] X. Qin, X. Tan, and S. Chen, “Tri-subject kinship verification: Understanding core of a family,” *CoRR arXiv:1501.02555*, 2015.
- [114] J. Liang, Q. Hu, C. Dang, and W. Zuo, “Weighted graph embedding-based metric learning for kinship verification,” *IEEE Trans. on Image Processing (TIP)*, vol. 28, no. 3, pp. 1149–1162, 2018.
- [115] H. Yan and J. Hu, “Video-based kinship verification using distance metric learning,” *Pattern Recognition*, 2018.
- [116] X. Wu, E. Boutellaa, X. Feng, and A. Hadid, “Kinship verification from faces: Methods, databases and challenges,” in *Conference on Signal Processing, Communications and Computing*. IEEE, 2016.
- [117] M. Georgopoulos, Y. Panagakis, and M. Pantic, “Modeling of facial aging and kinship: A survey,” *Image & Vision Computing*, 2018.
- [118] X. Qin, D. Liu, and D. Wang, “A literature survey on kinship verification through facial images,” *Neurocomputing*, 2019.
- [119] S. Xia, M. Shao, J. Luo, and Y. Fu, “Understanding kin relationships in a photo,” *IEEE Trans. on Multimedia*, vol. 14, no. 4, pp. 1046–1056, 2012.

## BIBLIOGRAPHY

- [120] J. Lu, J. Hu, X. Zhou, J. Zhou, M. Castrillón-Santana, J. Lorenzo-Navarro, L. Kou, Y. Shang, A. Bottino, and T. Figuieiredo Vieira, “Kinship verification in the wild: The first kinship verification competition,” in *IEEE International Joint Conference on Biometrics*, 2014.
- [121] M. B. Lopez, A. Hadid, E. Boutellaa, J. Goncalves, V. Kostakos, and S. Hosio, “Kinship verification from facial images and videos: human versus machine,” *Machine Vision and Applications*, 2018.
- [122] H. Dibeklioğlu, A. A. Salah, and T. Gevers, “Are you really smiling at me? spontaneous versus posed enjoyment smiles,” in *European Conference on Computer Vision (ECCV)*. Springer, 2012.
- [123] D. Hettachchi, N. van Berkel, S. Hosio, M. B. López, V. Kostakos, and J. Goncalves, “Augmenting automated kinship verification with targeted human input.” in *PACIS*, 2020, p. 141.
- [124] D. Aspandi, O. Martinez, and X. Binefa, “Heatmap-guided balanced deep convolution networks for family classification in the wild,” in *Conference on Automatic Face and Gesture Recognition*, 2019.
- [125] J. Lu, X. Zhou, Y.-P. Tan, Y. Shang, and J. Zhou, “Neighborhood repulsed metric learning for kinship verification,” *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, 2014.
- [126] H. Yan, J. Lu, and X. Zhou, “Prototype-based discriminative feature learning for kinship verification,” *IEEE Transactions on Cybernetics*, vol. 45, no. 11, pp. 2535–2545, 2014.
- [127] H. Yan, J. Lu, W. Deng, and X. Zhou, “Discriminative multimetric learning for kinship verification,” *IEEE Trans. on Information Forensics and Security*, vol. 9, no. 7, pp. 1169–1178, 2014.
- [128] X. Zhou, H. Yan, and Y. Shang, “Kinship verification from facial images by scalable similarity fusion,” *Neurocomputing*, 2016.
- [129] Y. Fang, Y. Y. S. Chen, H. Wang, and C. Shu, “Sparse similarity metric learning for kinship verification,” in *Visual Communications and Image Processing (VCIP)*. IEEE, 2016, pp. 1–4.
- [130] H. Liu and C. Zhu, “Status-aware projection metric learning for kinship verification,” in *IEEE International Conference on Multimedia and Expo (ICME)*, 2017, pp. 319–324.

## BIBLIOGRAPHY

- [131] X. Wu, X. Feng, E. Boutellaa, and A. Hadid, “Kinship verification using color features and extreme learning machine,” in *International Conference on Signal and Image Processing*. IEEE, 2018.
- [132] N. Kohli, M. Vatsa, R. Singh, A. Noore, and A. Majumdar, “Hierarchical representation learning for kinship verification,” *IEEE Trans. on Image Processing (TIP)*, vol. 26, no. 1, pp. 289–302, 2016.
- [133] J. Hu, J. Lu, L. Liu, and J. Zhou, “Multi-view geometric mean metric learning for kinship verification,” in *International Conference on Image Processing (ICIP)*. IEEE, 2019, pp. 1178–1182.
- [134] J. Lu, J. Hu, and Y.-P. Tan, “Discriminative deep metric learning for face and kinship verification,” *IEEE Trans. on Image Processing (TIP)*, vol. 26, no. 9, pp. 4269–4282, 2017.
- [135] X. Zhou, K. Jin, M. Xu, and G. Guo, “Learning deep compact similarity metric for kinship verification from face images,” *Information Fusion*, vol. 48, pp. 84–94, 2019.
- [136] O. Laiadi, A. Ouamane, A. Benakcha, A. Taleb-Ahmed, and A. Hadid, “Multi-view deep features for robust facial kinship verification,” *CoRR arXiv:2006.01315*, 2020.
- [137] X. Wu, E. Granger, T. H. Kinnunen, X. Feng, and A. Hadid, “Audio-visual kinship verification in the wild,” *IEEE International Conference on Biometrics (ICB)*, 2019.
- [138] S. Walsh and M. Kayser, “Predicting human appearance from dna for forensic investigations,” *Handbook of Forensic Genetics: Biodiversity and Heredity in Civil and Criminal Investigation*. New Jersey: World Scientific, pp. 415–448, 2016.
- [139] S. X. M. Shao and Y. Fu, “Genealogical face recognition based on ub kinface database,” in *IEEE CVPR Workshop on Biometrics*, 2011.
- [140] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, “Imagenet large scale visual recognition challenge,” *International Journal of Computer Vision (IJCV)*, 2015.
- [141] L. Wolf, “Deepface: Closing the gap to human-level performance in face verification,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.

## BIBLIOGRAPHY

- [142] F. Schroff, D. Kalenichenko, and J. Philbin, “Facenet: A unified embedding for face recognition and clustering,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 815–823.
- [143] K. Zhang, Y. Huang, C. Song, H. Wu, and L. Wang, “Kinship verification with deep convolutional neural networks,” in *Proceedings of the British Machine Vision Conference (BMVC)*. BMVA Press, September 2015, pp. 148.1–148.12.
- [144] C. Xiong, L. Liu, X. Zhao, S. Yan, and T.-K. Kim, “Convolutional fusion network for face verification in the wild,” 2015.
- [145] H. Liu and Y. Fu, “Clustering with partition level side information,” in *Proceedings of International Conference on Data Mining*, 2015.
- [146] C. W. Leong, R. Mihalcea, and S. Hassan, “Text mining for automatic image tagging,” in *Proceedings of the 23rd International Conference on Computational Linguistics: Posters*, ser. COLING ’10. PA, USA: Association for Computational Linguistics, 2010.
- [147] E. Law, B. Settles, and T. Mitchell, “Learning to tag using noisy labels,” in *Proc. ECML*, 2010, pp. 1–29.
- [148] D. Wang, S. C. Hoi, Y. He, and J. Zhu, “Mining weakly labeled web facial images for search-based face annotation,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 26, no. 1, pp. 166–179, 2014.
- [149] D. Yi, Z. Lei, S. Liao, and S. Z. Li, “Learning face representation from scratch,” *CoRR arXiv:1411.7923*, 2014.
- [150] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller, “Labeled faces in the wild: A database for studying face recognition in unconstrained environments,” UMass, Amherst, Tech. Rep., 2007.
- [151] D. E. King, “Dlib-ml: A machine learning toolkit,” *Journal of Machine Learning Research*, vol. 10, pp. 1755–1758, 2009.
- [152] J. R. Finkel, T. Grenager, and C. Manning, “Incorporating non-local information into information extraction systems by gibbs sampling,” in *Proceedings of the 43rd annual meeting on association for computational linguistics*. Association for Computational Linguistics, 2005.

## BIBLIOGRAPHY

- [153] H. Liu, Z. Tao, and Y. Fu, “Partition level constrained clustering,” *IEEE transactions on pattern analysis and machine intelligence*, 2017.
- [154] B. Mirkin, “Reinterpreting the category utility function,” *Machine Learning*, 2001.
- [155] J. Wu, H. Liu, H. Xiong, J. Cao, and J. Chen, “K-means-based consensus clustering: A unified view,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, no. 1, pp. 155–169, 2015.
- [156] T. Ahonen, A. Hadid, and M. Pietikainen, “Face description with local binary patterns: Application to face recognition,” *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, 2006.
- [157] N. Dalal and B. Triggs, “Histograms of oriented gradients for human detection,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 1, June 2005, pp. 886–893 vol. 1.
- [158] Y. Wen, K. Zhang, Z. Li, and Y. Qiao, “A discriminative feature learning approach for deep face recognition,” in *European Conference on Computer Vision (ECCV)*. Springer, 2016.
- [159] J. V. Davis, B. Kulis, P. Jain, S. Sra, and I. S. Dhillon, “Information-theoretic metric learning,” in *ICML*. ACM, 2007, pp. 209–216.
- [160] X. Niyogi, “Locality preserving projections,” in *Advances in Neural Information Processing Systems (NIPS)*. MIT, 2004.
- [161] K. Q. Weinberger and L. K. Saul, “Distance metric learning for large margin nearest neighbor classification,” *Journal of Machine Learning Research*, vol. 10, no. Feb, pp. 207–244, 2009.
- [162] Y. Peng, S. Wang, and B.-L. Lu, “Marginalized denoising autoencoder via graph regularization for domain adaptation,” in *Advances in Neural Information Processing Systems (NIPS)*. Springer, 2013.
- [163] Z. Ding, S. Suh, J.-J. Han, C. Choi, and Y. Fu, “Discriminative low-rank metric learning for face recognition,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2015.
- [164] D. G. Lowe, “Distinctive image features from scale-invariant keypoints,” *International Journal of Computer Vision (IJCV)*, vol. 60, 2004.

## BIBLIOGRAPHY

- [165] L. Wolf, T. Hassner, and I. Maoz, “Face recognition in unconstrained videos with matched background similarity,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2011.
- [166] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller, “Labeled faces in the wild: A database for studying face recognition in unconstrained environments,” Technical Report 07-49, University of Massachusetts, Amherst, Tech. Rep., 2007.
- [167] D. Pelleg and D. Baras, “K-means with large and noisy constraint sets,” in *Proceedings of European Conference on Machine Learning*, 2007, pp. 674–682.
- [168] M. F. Dal Martello and L. T. Maloney, “Where are kin recognition signals in the human face?” *Journal of Vision*, vol. 6, no. 12, 2006.
- [169] L. Best-Rowden, S. Bisht, J. C. Klontz, and A. K. Jain, “Unconstrained face recognition: Establishing baseline human performance via crowdsourcing,” in *International Journal of Computational Methods (IJCM)*. IEEE, 2014.
- [170] Q. Duan and L. Zhang, “Advnet: Adversarial contrastive residual net for 1 million kinship recognition,” in *RFIW at ACM MM*, 2017.
- [171] Y. Li, J. Zeng, J. Zhang, A. Dai, M. Kan, S. Shan, and X. Chen, “Kinnet: Fine-to-coarse deep metric learning for kinship verification,” in *RFIW at ACM MM*, 2017.
- [172] F. Wang, J. Cheng, W. Liu, and H. Liu, “Additive margin softmax for face verification,” *IEEE Signal Processing Letters*, vol. 25, no. 7, pp. 926–930, 2018.
- [173] S. Hörmann, M. Knoche, and G. Rigoll, “A multi-task comparator framework for kinship verification,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2020.
- [174] J. Yu, G. Xie, M. Li, and X. Hao, “Deep fusion siamese network for automatic kinship verification,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2020.
- [175] L. Zhipeng, Z. Zhiguang, X. Zhenyu, and C. Lixuan, “Challenge report-recognizing families in the wild data challenge,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2020.

## BIBLIOGRAPHY

- [176] A. Shadrikov, “Achieving better kinship recognition through better baseline,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2020.
- [177] H. V. Nguyen and L. Bai, “Cosine similarity metric learning for face verification,” in *Asian Conference on Computer Vision*. Springer, 2010, pp. 709–720.
- [178] C. Whitelam, E. Taborsky, A. Blanton, B. Maze, J. Adams, T. Miller, N. Kalka, A. K. Jain, J. A. Duncan, K. Allen *et al.*, “Iarpa janus benchmark-b face dataset,” in *Conference on Computer Vision and Pattern Recognition Workshop*, 2017.
- [179] J. Liang, J. Guo, S. Lao, and J. Li, “Using deep relational features to verify kinship,” in *CCF Chinese Conference on Computer Vision*. Springer, 2017, pp. 563–573.
- [180] X. Qin, D. Liu, and D. Wang, “Heterogeneous similarity learning for more practical kinship verification,” *Neural Processing Letters*, vol. 47, no. 3, pp. 1253–1269, 2018.
- [181] Z. Wei, M. Xu, L. Geng, H. Liu, and H. Yin, “Adversarial similarity metric learning for kinship verification,” *IEEE Access*, 2019.
- [182] S. Zhu, T. Yang, and C. Chen, “Visual explanation for deep metric learning,” *CoRR arXiv:1909.12977*, 2019.
- [183] M. Mukherjee and T. Meenpal, “Kinship verification using compound local binary pattern and local feature discriminant analysis,” in *2019 10th International Conference on Computing, Communication and Networking Technologies (ICCCNT)*. IEEE, 2019, pp. 1–7.
- [184] L. Zhang, Q. Duan, D. Zhang, W. Jia, and X. Wang, “Advkin: Adversarial convolutional network for kinship verification,” *IEEE Transactions on Cybernetics*, 2020.
- [185] O. Laiadi, A. Ouamane, A. Benakcha, A. Taleb-Ahmed, and A. Hadid, “Tensor cross-view quadratic discriminant analysis for kinship verification in the wild,” *Neurocomputing*, vol. 377, pp. 286–300, 2020.
- [186] W. Wang, S. You, S. Karaoglu, and T. Gevers, “Kinship identification through joint learning using kinship verification ensemble,” *CoRR arXiv:2004.06382*, 2020.
- [187] M. Wang, X. Shu, J. Feng, X. Wang, and J. Tang, “Deep multi-person kinship matching and recognition for family photos,” *Pattern Recognition*, p. 107342, 03 2020.

## BIBLIOGRAPHY

- [188] A. Grouver, P. Shivakumara, M. A. Kaljahi, B. Chetty, U. Pal, T. Lu, and G. H. Kumar, “A spatial density and phase angle based correlation for multi-type family photo identification,” in *Asian Conference on Pattern Recognition*. Springer, 2019, pp. 76–89.
- [189] F. Crispim, T. Vieira, and B. Lima, “Verifying kinship from rgb-d face data,” in *International Conference on Advanced Concepts for Intelligent Vision Systems*. Springer, 2020, pp. 215–226.
- [190] Y. B. Kumar Ravi, C. K. Narayananappa, and P. Dayanandab, “Weighted full binary tree-sliced binary pattern: An rgb-d image descriptor,” *Heliyon*, vol. 6, 2020.
- [191] S. Mahpod and Y. Keller, “Kinship verification using multiview hybrid distance learning,” *Computer Vision and Image Understanding*, vol. 167, pp. 28–36, 2018.
- [192] J. Hu, J. Lu, L. Liu, and J. Zhou, “Multi-view geometric mean metric learning for kinship verification,” in *International Conference on Image Processing (ICIP)*. IEEE, 2019, pp. 1178–1182.
- [193] J. Hu, J. Lu, Y.-P. Tan, J. Yuan, and J. Zhou, “Local large-margin multi-metric learning for face and kinship verification,” *Transactions on Circuits and Systems for Video Technology*, 2017.
- [194] M. A. Kaljahi, P. Shivakumara, T. Hu, H. A. Jalab, R. W. Ibrahim, M. Blumenstein, T. Lu, and M. N. B. Ayub, “A geometric and fractional entropy-based method for family photo classification,” *Expert Systems with Applications: X*, 2019.
- [195] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman, “Vggface2: A dataset for recognising faces across pose and age,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*. IEEE, 2018.
- [196] Y. Guo, L. Zhang, Y. Hu, X. He, and J. Gao, “Ms-celeb-1m: A dataset and benchmark for large-scale face recognition,” in *European Conference on Computer Vision (ECCV)*. Springer, 2016.
- [197] J. Bromley, I. Guyon, Y. LeCun, E. Säckinger, and R. Shah, “Signature verification using a “siamese” time delay neural network,” in *Advances in Neural Information Processing Systems (NIPS)*, 1994.

## BIBLIOGRAPHY

- [198] Z. Zhang, Y. Song, and H. Qi, “Age progression/regression by conditional adversarial autoencoder,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2017.
- [199] J. Cheng, Y. Li, J. Wang, L. Yu, and S. Wang, “Exploiting effective facial patches for robust gender recognition,” *Tsinghua Science and Technology*, vol. 24, no. 3, pp. 333–345, 2019.
- [200] M. Wang, W. Deng, J. Hu, J. Peng, X. Tao, and Y. Huang, “Racial faces in-the-wild: Reducing racial bias by deep unsupervised domain adaptation,” *CoRR arXiv:1812.00194*, 2018.
- [201] Q. Zhang, L. T. Yang, Z. Chen, and P. Li, “A survey on deep learning for big data,” *Information Fusion*, vol. 42, pp. 146–157, 2018.
- [202] T.-D. H. Nguyen, H.-N. Nguyen, and H. Dao, “Recognizing families through images,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2020.
- [203] J. Yu, G. Xie, M. Li, and X. Hao, “Retrieval of family members using siamese neural network,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2020.
- [204] K. Zhang12, Y. Huang, C. Song, H. Wu, and L. Wang, “Kinship verification with deep convolutional neural networks,” *British Machine Vision Conference (BMVC)*, 2015.
- [205] X. Qin, D. Liu, and D. Wang, “Social relationships classification using social contextual features and svdd-based metric learning,” *Applied Soft Computing*, vol. 77, pp. 344–355, 2019.
- [206] Q. Duan, L. Zhang, and W. Zuo, “From face recognition to kinship verification: An adaptation approach,” in *Conference on Computer Vision and Pattern Recognition Workshop*, 2017.
- [207] H. Zhang, X. Wang, and C.-C. J. Kuo, “Deep kinship verification via appearance-shape joint prediction and adaptation-based approach,” in *International Conference on Image Processing (ICIP)*, 2019.
- [208] A. Chergui, S. Ouchtati, S. Mavromatis, S. Bekhouche, M. Lashab, and J. Sequeira, “Kinship verification through facial images using cnn-based features,” *Traitement du Signal*, vol. 37, no. 1, pp. 1–8, 2020.
- [209] X. Chu, Z. Cheng, J. Zhao, X. Huang, and W.-L. Wu, “Kinship verification with an optimal cnn,” 2017, unpublished paper.

## BIBLIOGRAPHY

- [210] N. Kohli, D. Yadav, M. Vatsa, R. Singh, and A. Noore, “Deep face-representation learning for kinship verification,” in *Deep Learning in Biometrics*. CRC Press, 2018, pp. 127–152.
- [211] Y. Sun, J. Li, Y. Wei, and H. Yan, “Video-based parent-child relationship prediction,” in *2018 IEEE Visual Communications and Image Processing (VCIP)*. IEEE, 2018, pp. 1–4.
- [212] J. Lv, B. Wu, Y. Zhang, and Y. Xiao, “Attentive sequences recurrent network for social relation recognition from video,” *IEICE Transaction on Information and Systems*, 2019.
- [213] S. Gong, Y. Shi, and A. K. Jain, “Video face recognition: Component-wise feature aggregation net,” *CoRR arXiv:1902.07327*, 2019.
- [214] W. Li, Y. Zhang, K. Lv, J. Lu, J. Feng, and J. Zhou, “Graph-based kinship reasoning network,” in *IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2020, pp. 1–6.
- [215] E. Dahan and Y. Keller, “Selfkin: self adjusted deep model for kinship verification,” *CoRR arXiv:1809.08493*, 2018.
- [216] O. Laiadi, A. Ouamane, A. Benakcha, A. Taleb-Ahmed, and A. Hadid, “Kinship verification based deep and tensor features through extreme learning machine,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*. IEEE, 2019, pp. 1–4.
- [217] A. Nandy and S. S. Mondal, “Kinship verification using deep siamese convolutional neural network,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*. IEEE, 2019.
- [218] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, “SqueezeNet: Alexnet-level accuracy with 50x fewer parameters and 0.5 mb model size,” *CoRR arXiv:1602.07360*, 2016.
- [219] S. Ozkan and A. Ozkan, “Kinshipgan: Synthesizing of kinship faces from family photos by regularizing a deep face network,” in *International Conference on Image Processing (ICIP)*, 2018.
- [220] I. Ö. Ertugrul and H. Dibeklioglu, “What will your future child look like? modeling and synthesis of hereditary patterns of facial dynamics,” in *Conference on Automatic Face and Gesture Recognition*, 2017.

## BIBLIOGRAPHY

- [221] M. Merler, N. Ratha, R. S. Feris, and J. R. Smith, “Diversity in faces,” *CoRR arXiv:1901.10436*, 2019.
- [222] M.-K. Hu, “Visual pattern recognition by moment invariants,” *IRE transactions on information theory*, 1962.
- [223] A. Ephrat, I. Mosseri, O. Lang, T. Dekel, K. Wilson, A. Hassidim, W. T. Freeman, and M. Rubinstein, “Looking to listen at the cocktail party: A speaker-independent audio-visual model for speech separation,” *CoRR arXiv:1804.03619*, 2018.
- [224] A. Nagrani, J. S. Chung, and A. Zisserman, “Voxceleb: a large-scale speaker identification dataset,” *CoRR arXiv:1706.08612*, 2017.
- [225] J. S. Chung, A. Nagrani, and A. Zisserman, “Voxceleb2: Deep speaker recognition,” *CoRR arXiv:1806.05622*, 2018.
- [226] A. Nagrani, S. Albanie, and A. Zisserman, “Seeing voices and hearing faces: Cross-modal biometric matching,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [227] S. Albanie, A. Nagrani, A. Vedaldi, and A. Zisserman, “Emotion recognition in speech using cross-modal transfer in the wild,” in *ACM Conference on Multimedia (ACMMM)*, 2018.
- [228] M. Hao, W.-H. Cao, Z.-T. Liu, M. Wu, and P. Xiao, “Visual-audio emotion recognition based on multi-task and ensemble learning with multiple features,” *Neurocomputing*, 2020.
- [229] O. Wiles, A. Koepke, and A. Zisserman, “X2face: A network for controlling face generation by using images, audio, and pose codes,” in *European Conference on Computer Vision (ECCV)*, 2018.
- [230] ——, “Self-supervised learning of a facial attribute embedding from video,” in *British Machine Vision Conference (BMVC)*, 2018.
- [231] X. Song, H. Chen, Q. Wang, Y. Chen, M. Tian, and H. Tang, “A review of audio-visual fusion with machine learning,” in *Journal of Physics: Conference Series*, vol. 1237, no. 2. IOP Publishing, 2019, p. 022144.
- [232] S. Petridis, Y. Wang, Z. Li, and M. Pantic, “End-to-end audiovisual fusion with lstms,” *CoRR arXiv:1709.04343*, 2017.

## BIBLIOGRAPHY

- [233] D. Zuo and P. Mok, “Formant dynamics of bilingual identical twins in non-contemporaneous speech,” in *Proceedings of Australasian International Conference on SST. Sydney, Australia*, 2012.
- [234] P. G. Hepper, “Long-term retention of kinship recognition established during infancy in the domestic dog,” *Behavioural processes*, vol. 33, no. 1-2, pp. 3–14, 1994.
- [235] P. Poindron, A. Terrazas, M. Oca, N. Serafín, and H. Hernandez, “Sensory and physiological determinants of maternal behavior in the goat (*capra hircus*),” *Hormones and behavior*, vol. 52, pp. 99–105, 07 2007.
- [236] P. Poindron, F. Lévy, and M. Keller, “Maternal responsiveness and maternal selectivity in domestic sheep and goats: the two facets of maternal attachment,” *Developmental Psychobiology: The Journal of the International Society for Developmental Psychobiology*, vol. 49, no. 1, pp. 54–70, 2007.
- [237] X. Zhou, J. Hu, J. Lu, Y. Shang, and Y. Guan, “Kinship verification from facial images under uncontrolled conditions,” in *Proceedings of the 19th ACM international conference on Multimedia*. ACM, 2011.
- [238] L. Zhang, K. Ma, H. Nejati, L. Foo, T. Sim, and D. Guo, “A talking profile to distinguish identical twins,” *Image and Vision Computing*, 2014.
- [239] M. Georgopoulos, Y. Panagakis, and M. Pantic, “Investigating bias in deep face analysis: The kanface dataset and empirical study,” *CoRR arXiv:2005.07302*, 2020.
- [240] F. S. Ghatas and E. E. Hemayed, “Gankin: generating kin faces using disentangled gan,” *SN Applied Sciences*, vol. 2, no. 2, pp. 1–10, 2020.
- [241] C. Kumar, R. Ryan, and M. Shao, “Adversary for social good: Protecting familial privacy through joint adversarial attacks,” in *Conference on Artificial Intelligence (AAAI)*, 2020.
- [242] B. Maze, J. Adams, J. A. Duncan, N. Kalka, T. Miller, C. Otto, A. K. Jain, W. T. Niggel, J. Anderson, J. Cheney *et al.*, “Iarpa janus benchmark-c: Face dataset and protocol,” in *IEEE International Conference on Biometrics (ICB)*. IEEE, 2018.
- [243] Y. Li, J. Zeng, J. Zhang, A. Dai, M. Kan, S. Shan, and X. Chen, “Kinnet: Fine-to-coarse deep metric learning for kinship verification,” in *RFIW at ACM MM*, 2017.

## BIBLIOGRAPHY

- [244] Q. Duan and L. Zhang, “Advnet: Adversarial contrastive residual net for 1 million kinship recognition,” in *RFIW at ACM MM*, 2017.
- [245] T. Afouras, J. S. Chung, and A. Zisserman, “The conversation: Deep audio-visual speech enhancement,” *CoRR arXiv:1804.04121*, 2018.
- [246] J. S. Chung and A. Zisserman, “Lip reading in profile,” in *British Machine Vision Conference (BMVC)*, 2017.
- [247] T. Afouras, J. S. Chung, and A. Zisserman, “Deep lip reading: a comparison of models and an online application,” in *INTERSPEECH*, 2018.
- [248] J. S. Chung and A. Zisserman, “Lip reading in the wild,” in *Asian Conference on Computer Vision*, 2016.
- [249] J. S. Chung, A. Senior, O. Vinyals, and A. Zisserman, “Lip reading sentences in the wild,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [250] J. S. Chung, A. Jamaludin, and A. Zisserman, “You said that?” in *British Machine Vision Conference (BMVC)*, 2017.
- [251] O. Wiles, A. Sophia Koepke, and A. Zisserman, “X2face: A network for controlling face generation using images, audio, and pose codes,” in *European Conference on Computer Vision (ECCV)*, 2018, pp. 670–686.
- [252] S. Albanie, A. Nagrani, A. Vedaldi, and A. Zisserman, “Emotion recognition in speech using cross-modal transfer in the wild,” in *ACM Conference on Multimedia (ACMMM)*, 2018.
- [253] I. U. Haq, K. Muhammad, T. Hussain, S. Kwon, M. Sodanil, S. W. Baik, and M. Y. Lee, “Movie scene segmentation using object detection and set theory,” *International Journal of Distributed Sensor Networks*, vol. 15, no. 6, p. 1550147719845277, 2019.
- [254] E. Sánchez-Nielsen, F. Chávez-Gutiérrez, J. Lorenzo-Navarro, and M. Castrillón-Santana, “A multimedia system to produce and deliver video fragments on demand on parliamentary websites,” *Multimedia Tools and Applications*, vol. 76, no. 5, pp. 6281–6307, 2017.
- [255] F. Wang, X. Xiang, J. Cheng, and A. L. Yuille, “Normface: L2 hypersphere embedding for face verification,” in *ACM Conference on Multimedia (ACMMM)*, 2017, pp. 1041–1049.

## BIBLIOGRAPHY

- [256] J. S. Chung, J. Huh, S. Mun, M. Lee, H. S. Heo, S. Choe, C. Ham, S. Jung, B.-J. Lee, and I. Han, “In defence of metric learning for speaker recognition,” *CoRR arXiv:2003.11982*, 2020.
- [257] H. Zhu, M. Luo, R. Wang, A. Zheng, and R. He, “Deep audio-visual learning: A survey,” *CoRR arXiv:2001.04758*, 2020.
- [258] J. S. Chung and A. Zisserman, “Out of time: automated lip sync in the wild,” in *Workshop on Multi-view Lip-reading, ACCV*, 2016.
- [259] Y. Li, M. Murias, S. Major, G. Dawson, K. Dzirasa, L. Carin, and D. E. Carlson, “Targeting eeg/lfp synchrony with neural nets,” in *Advances in Neural Information Processing Systems (NIPS)*, 2017, pp. 4621–4631.
- [260] B. DeCann and A. Ross, “Relating roc and cmc curves via the biometric menagerie,” in *2013 IEEE Sixth International Conference on Biometrics: Theory, Applications and Systems (BTAS)*. IEEE, 2013, pp. 1–8.
- [261] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *CoRR arXiv:1412.6980*, 2014.
- [262] J. Snell, K. Swersky, and R. Zemel, “Prototypical networks for few-shot learning,” in *Advances in Neural Information Processing Systems (NIPS)*, 2017, pp. 4077–4087.
- [263] L. Wan, Q. Wang, A. Papir, and I. L. Moreno, “Generalized end-to-end loss for speaker verification,” in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 4879–4883.
- [264] J. Wang, F. Zhou, S. Wen, X. Liu, and Y. Lin, “Deep metric learning with angular loss,” in *IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2593–2601.
- [265] B. Schölkopf, A. J. Smola, F. Bach *et al.*, *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT, 2002.
- [266] D. E. King, “Dlib-ml: A machine learning toolkit,” *The Journal of Machine Learning Research*, vol. 10, pp. 1755–1758, 2009.

## BIBLIOGRAPHY

- [267] L. Xiong, J. Karlekar, J. Zhao, Y. Cheng, Y. Xu, J. Feng, S. Pranata, and S. Shen, “A good practice towards top performance of face recognition: Transferred deep feature fusion,” *CoRR arXiv:1704.00438*, 2017.
- [268] M. B. López, E. Boutellaa, and A. Hadid, “Comments on the “kinship face in the wild” data sets,” *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, 2016.
- [269] M. F. Dal Martello, L. M. DeBruine, and L. T. Maloney, “Allocentric kin recognition is not affected by facial inversion,” *Journal of vision*, vol. 15, no. 13, pp. 5–5, 2015.
- [270] G. Kaminski, F. Ravary, C. Graff, and E. Gentaz, “Firstborns’ disadvantage in kinship detection,” *Psychological science*, 2010.
- [271] X. Wu, E. Boutellaa, M. B. López, X. Feng, and A. Hadid, “On the usefulness of color for kinship verification from face images,” in *IEEE workshop on information forensics and security*, 2016.
- [272] Z. McNatt, N. G. Boothby, H. Al-Shannaq, H. Chandler, P. Freels, A. S. Mahmoud, N. Majdalani, and L. Zebib, “Impact of separation on refugee families: Syrian refugees in jordan,” 2018.
- [273] “The ethics of catching criminals using their family’s dna,” *Nature*, 2018.
- [274] F. Taherkhani, N. M. Nasrabadi, and J. Dawson, “A deep face identification network enhanced by facial attributes prediction,” in *Conference on Computer Vision and Pattern Recognition Workshop*, 2018.
- [275] F. B. Naini and J. P. Moss, “Three-dimensional assessment of the relative contribution of genetics and environment to various facial parameters with the twin method,” *American Journal of Orthodontics and Dentofacial Orthopedics*, vol. 126, no. 6, pp. 655–665, 2004.
- [276] C. D. Froud, V. Bruce, H. Y. Chang, Y. Plenderleith, A. H. McIntyre, and P. J. Hancock, “Predict your child: a system to suggest the facial appearance of children,” *Journal of Multimedia*, 2008.
- [277] A. Acien, A. Morales, R. Vera-Rodriguez, I. Bartolome, and J. Fierrez, “Measuring the gender and ethnicity bias in deep models for face recognition,” in *Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications*, R. Vera-Rodriguez, J. Fierrez, and A. Morales, Eds. Springer International Publishing, 2019.

## BIBLIOGRAPHY

- [278] L. Anne Hendricks, K. Burns, K. Saenko, T. Darrell, and A. Rohrbach, “Women also snowboard: Overcoming bias in captioning models,” in *European Conference on Computer Vision (ECCV)*, 2018.
- [279] M. Wang, W. Deng, J. Hu, J. Peng, X. Tao, and Y. Huang, “Race faces in-the-wild: Reduce bias by deep unsupervised domain adaptation,” *CoRR arXiv:1812.00194*, 2019.
- [280] S. Nagpal, M. Singh, R. Singh, M. Vatsa, and N. Ratha, “Deep learning for face recognition: Pride or prejudiced?” *CoRR arXiv:1904.01219*, 2019.
- [281] J. Snow, “Amazon’s face recognition falsely matched 28 members of congress with mugshots,” 2018. [Online]. Available: <https://www.aclu.org/blog/privacy-technology/surveillance-technologies/amazons-face-recognition-falsely-matched-28>
- [282] C. Lazo, “Towards engineering ai software for fairness: A framework to help design fair, accountable and transparent algorithmic decision-making systems,” *TU Delft Library*, 2020.
- [283] P. Drozdowski, C. Rathgeb, A. Dantcheva, N. Damer, and C. Busch, “Demographic bias in biometrics: A survey on an emerging challenge,” *IEEE Transactions on Technology and Society*, 2020.
- [284] S. Windmann and T. Krüger, “Subconscious detection of threat as reflected by an enhanced response bias,” *Consciousness and cognition*, 1998.
- [285] T. Tommasi, N. Patricia, B. Caputo, and T. Tuytelaars, “A deeper look at dataset bias,” in *Domain adaptation in computer vision applications*. Springer, 2017.
- [286] A. Amini, A. P. Soleimany, W. Schwarting, S. N. Bhatia, and D. Rus, “Uncovering and mitigating algorithmic bias through learned latent structure,” in *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, 2019, pp. 289–295.
- [287] B. Kim, H. Kim, K. Kim, S. Kim, and J. Kim, “Learning not to learn: Training deep neural networks with biased data,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [288] P. Stock and M. Cisse, “Convnets and imagenet beyond accuracy: Understanding mistakes and uncovering biases,” in *European Conference on Computer Vision (ECCV)*, 2018.

## BIBLIOGRAPHY

- [289] G. Balakrishnan, Y. Xiong, W. Xia, and P. Perona, “Towards causal benchmarking of bias in face analysis algorithms,” *CoRR arXiv:2007.06570*, 2020.
- [290] M. Georgopoulos, J. Oldfield, M. A. Nicolaou, Y. Panagakis, and M. Pantic, “Enhancing facial data diversity with style-based face aging,” in *Conference on Computer Vision and Pattern Recognition Workshop*, 2020.
- [291] P. Terhörst, J. N. Kolf, N. Damer, F. Kirchbuchner, and A. Kuijper, “Post-comparison mitigation of demographic bias in face recognition using fair score normalization,” *CoRR arXiv:2002.03592*, 2020.
- [292] I. Serna, A. Peña, A. Morales, and J. Fierrez, “Insidebias: Measuring bias in deep networks and application to face gender biometrics,” *CoRR arXiv:2004.06592*, 2020.
- [293] V. Albiero, K. KS, K. Vangara, K. Zhang, M. C. King, and K. W. Bowyer, “Analysis of gender inequality in face recognition accuracy,” in *Conference on Computer Vision and Pattern Recognition Workshop*, 2020, pp. 81–89.
- [294] A. Das, A. Dantcheva, and F. Bremond, “Mitigating bias in gender, age and ethnicity classification: a multi-task convolution neural network approach,” in *European Conference on Computer Vision (ECCV)*, 2018.
- [295] N. Srinivas, K. Ricanek, D. Michalski, D. S. Bolme, and M. King, “Face recognition algorithm bias: Performance differences on images of children and adults,” in *Conference on Computer Vision and Pattern Recognition Workshop*, 2019.
- [296] A. Acien, A. Morales, R. Vera-Rodriguez, I. Bartolome, and J. Fierrez, “Measuring the gender and ethnicity bias in deep models for face recognition,” in *Iberoamerican Congress on Pattern Recognition*. Springer, 2018.
- [297] S. Gong, X. Liu, and A. K. Jain, “Debiasface: De-biasing face recognition,” *CoRR arXiv:1911.08080*, 2019.
- [298] A. V. Savchenko, “Efficient facial representations for age, gender and identity recognition in organizing photo albums using multi-output convnet,” *PeerJ Computer Science*, vol. 5, p. e197, 2019.

## BIBLIOGRAPHY

- [299] S. Nagpal, M. Singh, R. Singh, and M. Vatsa, “Attribute aware filter-drop for bias invariant classification,” in *Conference on Computer Vision and Pattern Recognition Workshop*, June 2020.
- [300] T. Sixta, J. Junior, C. Jacques, P. Buch-Cardona, E. Vazquez, and S. Escalera, “Fairface challenge at eccv 2020: Analyzing bias in face recognition,” *CoRR arXiv:2009.07838*, 2020.
- [301] A. Peña, I. Serna, A. Morales, and J. Fierrez, “Bias in multimodal ai: Testbed for fair automatic recruitment,” *CoRR arXiv:2004.07173*, 2020.
- [302] ———, “Faircvtest demo: Understanding bias in multimodal learning with testbed in fair automatic recruitment,” *CoRR arXiv:2009.07025*, 2020.
- [303] P. Terhörst, N. Damer, F. Kirchbuchner, and A. Kuijper, “Suppressing gender and age in face templates using incremental variable elimination,” in *2019 International Conference on Biometrics (ICB)*, 2019, pp. 1–8.
- [304] B. Sadeghi and V. N. Boddeti, “Imparting fairness to pre-trained biased representations,” in *Conference on Computer Vision and Pattern Recognition Workshop*, June 2020.
- [305] M. Wang and W. Deng, “Mitigating bias in face recognition using skewness-aware reinforcement learning,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [306] P.-M. Law, S. Malik, F. Du, and M. Sinha, “Designing tools for semi-automated detection of machine learning biases: An interview study,” *CoRR arXiv:2003.07680*, 2020.
- [307] X. Yin, X. Yu, K. Sohn, X. Liu, and M. Chandraker, “Feature transfer learning for face recognition with under-represented data,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [308] C. Huang, X. Chen, P. Kairouz, L. Sankar, and R. Rajagopal, “Generative adversarial models for learning private and fair representations,” 2018.
- [309] Y. Wang, Y. Feng, H. Liao, J. Luo, and X. Xu, “Do they all look the same? deciphering chinese, japanese and koreans by fine-grained deep learning,” in *IEEE Conference on Multimedia Information Processing and Retrieval*, 2018.

## BIBLIOGRAPHY

- [310] V. Muthukumar, T. Pedapati, N. Ratha, P. Sattigeri, C.-W. Wu, B. Kingsbury, A. Kumar, S. Thomas, A. Mojsilovic, and K. Varshney, “Understanding unequal gender classification from face images,” *CoRR arXiv:1812.00099*, 2018.
- [311] A. Das, A. Dantcheva, and F. Bremond, “Mitigating bias in gender, age and ethnicity classification: a multi-task convolution neural network approach,” *European Conference on Computer Vision (ECCV)*, 2018.
- [312] I. Hupont and C. Fernandez, “Demogpairs: Quantifying the impact of demographic imbalance in deep face recognition,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2019.
- [313] E. López-López, X. M. Pardo, C. V. Regueiro, R. Iglesias, and F. E. Casado, “Dataset bias exposed in face verification,” *IET Biometrics*, 2019.
- [314] A. Das, A. Dantcheva, and F. Bremond, “Nature and nurture in own-race face processing,” in *Psychological science*, 2006.
- [315] C. A. Meissner and J. C. Brigham, “Thirty years of investigating the own-race bias in memory for faces: A meta-analytic review,” in *Psychology, Public Policy, and Law*, 2001.
- [316] M. E. Nicholls, O. Churches, and T. Loetscher, “Perception of an ambiguous figure is affected by own-age social biases,” in *Scientific reports*, 2018.
- [317] C. Drummond, R. C. Holte *et al.*, “C4. 5, class imbalance, and cost sensitivity: why under-sampling beats over-sampling,” in *Workshop on learning from imbalanced datasets II*. Citeseer, 2003.
- [318] M. Oquab, L. Bottou, I. Laptev, and J. Sivic, “Learning and transferring mid-level image representations using convolutional neural networks,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014, pp. 1717–1724.
- [319] E. M. Rudd, M. Günther, and T. E. Boult, “Moon: A mixed objective optimization network for the recognition of facial attributes,” in *European Conference on Computer Vision (ECCV)*. Springer, 2016, pp. 19–35.
- [320] C. Huang, Y. Li, C. L. Chen, and X. Tang, “Deep imbalanced learning for face recognition and attribute prediction,” *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, 2019.

## BIBLIOGRAPHY

- [321] H. He and E. A. Garcia, “Learning from imbalanced data,” *IEEE Transactions on knowledge and data engineering*, 2009.
- [322] H. He and Y. Ma, *Imbalanced learning: foundations, algorithms, and applications*. John Wiley & Sons, 2013.
- [323] B. Krawczyk, “Learning from imbalanced data: open challenges and future directions,” *Progress in Artificial Intelligence*, 2016.
- [324] T. Xu, J. White, S. Kalkan, and H. Gunes, “Investigating bias and fairness in facial expression recognition,” *CoRR arXiv:2007.10075*, 2020.
- [325] V. Albiero, K. Zhang, and K. W. Bowyer, “How does gender balance in training data affect face recognition accuracy?” *CoRR arXiv:2002.02934*, 2020.
- [326] I. Hupont and C. Fernández, “Demogpairs: Quantifying the impact of demographic imbalance in deep face recognition,” in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*. IEEE, 2019, pp. 1–7.
- [327] K. Kärkkäinen and J. Joo, “Fairface: Face attribute dataset for balanced race, gender, and age,” *CoRR arXiv:1908.04913*, 2019.
- [328] Z. Ding, S. Li, M. Shao, and Y. Fu, “Graph adaptive knowledge transfer for unsupervised domain adaptation,” in *ECCV*, 2018.
- [329] X. Peng, B. Usman, N. Kaushik, J. Hoffman, D. Wang, and K. Saenko, “Visda: The visual domain adaptation challenge,” *CoRR arXiv:1710.06924*, 2017.
- [330] K. Saito, D. Kim, S. Sclaroff, T. Darrell, and K. Saenko, “Semi-supervised domain adaptation via minimax entropy,” *CoRR arXiv:1904.06487*, 2019.
- [331] R. Shu, H. H. Bui, H. Narui, and S. Ermon, “A dirt-t approach to unsupervised domain adaptation,” *CoRR arXiv:1802.08735*, 2018.
- [332] B. Sun and K. Saenko, “Subspace distribution alignment for unsupervised domain adaptation.” in *BMVC*, 2015.
- [333] M. Long, J. Wang, G. Ding, J. Sun, and P. S. Yu, “Transfer feature learning with joint distribution adaptation,” in *ICCV*, 2013.

## BIBLIOGRAPHY

- [334] B. Gong, Y. Shi, F. Sha, and K. Grauman, “Geodesic flow kernel for unsupervised domain adaptation,” in *CVPR*, 2012.
- [335] R. Gopalan, R. Li, and R. Chellappa, “Domain adaptation for object recognition: An unsupervised approach,” in *ICCV*, 2011.
- [336] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, “Domain-adversarial training of neural networks,” *JMLR*, vol. 17, no. 1, 2016.
- [337] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, “Adversarial discriminative domain adaptation,” in *CVPR*, 2017.
- [338] B. Güler, A. S. Avestimehr, and A. Ortega, “Privacy-aware distributed graph-based semi-supervised learning,” in *International Workshop on Machine Learning for Signal Processing*. IEEE, 2019.
- [339] K. W. Bowyer, “Face recognition technology: security versus privacy,” *IEEE Technology and society magazine*, vol. 23, no. 1, pp. 9–19, 2004.
- [340] P. Dhar, J. Gleason, H. Souri, C. D. Castillo, and R. Chellappa, “An adversarial learning algorithm for mitigating gender bias in face recognition,” *CoRR arXiv:2006.07845*, 2020.
- [341] V. Mirjalili, S. Raschka, and A. Ross, “Gender privacy: An ensemble of semi adversarial networks for confounding arbitrary gender classifiers,” in *2018 IEEE 9th International Conference on Biometrics Theory, Applications and Systems (BTAS)*, 2018, pp. 1–10.
- [342] A. Othman and A. Ross, “Privacy of facial soft biometrics: Suppressing gender but retaining identity,” in *Computer Vision - ECCV 2014 Workshops*, L. Agapito, M. M. Bronstein, and C. Rother, Eds. Cham: Springer International Publishing, 2015, pp. 682–696.
- [343] S. Guo, T. Xiang, and X. Li, “Towards efficient privacy-preserving face recognition in the cloud,” *Signal Processing*, vol. 164, pp. 320 – 328, 2019. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0165168419302324>
- [344] Z. Ma, Y. Liu, X. Liu, J. Ma, and K. Ren, “Lightweight privacy-preserving ensemble classification for face recognition,” *IEEE Internet of Things Journal*, vol. 6, no. 3, pp. 5778–5790, 2019.

## BIBLIOGRAPHY

- [345] A. Ambekar, C. Ward, J. Mohammed, S. Male, and S. Skiena, “Name-ethnicity classification from open sources,” in *Proceedings of the SIGKDD conference on Knowledge Discovery and Data Mining*, 2009.
- [346] S. Fu, H. He, and Z.-G. Hou, “Learning race from face: A survey,” *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, 2014.
- [347] G. Levi and T. Hassner, “Age and gender classification using convolutional neural networks,” in *Conference on Computer Vision and Pattern Recognition Workshop*, 2015, pp. 34–42.
- [348] R. England, “Facial recognition misidentified 26 california lawmakers as criminals,” 2019.
- [349] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: a simple way to prevent neural networks from overfitting,” *The journal of machine learning research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [350] L. A. Jeni, J. F. Cohn, and F. De La Torre, “Facing imbalanced data—recommendations for the use of performance metrics,” in *Humaine association conference on affective computing and intelligent interaction*. IEEE, 2013, pp. 245–251.
- [351] T. Gebru, J. Morgenstern, B. Vecchione, J. W. Vaughan, H. Wallach, H. Daumeé III, and K. Crawford, “Datasheets for datasets,” *CoRR arXiv:1803.09010*, 2018.
- [352] M. Arnold, R. K. Bellamy, M. Hind, S. Houde, S. Mehta, A. Mojsilović, R. Nair, K. N. Ramamurthy, A. Olteanu, D. Piorkowski *et al.*, “Factsheets: Increasing trust in ai services through supplier’s declarations of conformity,” *IBM Journal of R&D*, 2019.
- [353] M. Mitchell, S. Wu, A. Zaldivar, P. Barnes, L. Vasserman, B. Hutchinson, E. Spitzer, I. D. Raji, and T. Gebru, “Model cards for model reporting,” *CoRR arXiv:1810.03993*, 2018.
- [354] W. J. Goode, “World revolution and family patterns.” 1963.
- [355] R. Zhang, F. Nie, X. Li, and X. Wei, “Feature selection with multi-view data: A survey,” *Information Fusion*, 2019.
- [356] S. Brédart and R. M. French, “Do babies resemble their fathers more than their mothers? a failure to replicate christenfeld and hill (1995),” *Evolution and Human Behavior*, 1999.
- [357] D. Kollias, A. Schulc, E. Hajiyev, and S. Zafeiriou, “Analyze affective behavior in abaw 2020 competition,” *CoRR arXiv:2001.11409*, 2020.