

Improving deep facial phenotyping for ultra-rare disorder verification using model ensembles

Alexander Hustinx, Fabio Hellmann, Ömer Sümer, Behnam Javanmardi, Elisabeth André, Peter Krawitz, Tzung-Chien Hsieh

Angaben zur Veröffentlichung / Publication details:

Hustinx, Alexander, Fabio Hellmann, Ömer Sümer, Behnam Javanmardi, Elisabeth André, Peter Krawitz, and Tzung-Chien Hsieh. 2023. "Improving deep facial phenotyping for ultra-rare disorder verification using model ensembles." In 2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), Jan. 2-7, 2023, Waikoloa, HI, USA, edited by David Crandall, Boqing Gong, Yong Jae Lee, Richard Souvenir, Stella Yu, Tamara Berg, and Ryan Farrell, 5007-17. Piscataway, NJ: IEEE. <https://doi.org/10.1109/WACV56688.2023.00499>.



Improving Deep Facial Phenotyping for Ultra-rare Disorder Verification Using Model Ensembles

Alexander Hustinx¹, Fabio Hellmann², Ömer Sümer², Behnam Javanmardi¹,
Elisabeth André², Peter Krawitz¹, Tzung-Chien Hsieh^{1*}

¹ Institute for Genomic Statistics and Bioinformatics, University Hospital Bonn, University of Bonn

² Chair for Human-Centered Artificial Intelligence, University of Augsburg

{ahustinx, b-jav, pkrawitz, thsieh}@uni-bonn.de

{oemer.suemer, fabio.hellmann, andre}@informatik.uni-augsburg.de

Abstract

Rare genetic disorders affect more than 6% of the global population. Reaching a diagnosis is challenging because rare disorders are very diverse. Many disorders have recognizable facial features that are hints for clinicians to diagnose patients. Previous work, such as GestaltMatcher, utilized representation vectors produced by a DCNN similar to AlexNet to match patients in high-dimensional feature space to support “unseen” ultra-rare disorders. However, the architecture and dataset used for transfer learning in GestaltMatcher have become outdated. Moreover, a way to train the model for generating better representation vectors for unseen ultra-rare disorders has not yet been studied. Because of the overall scarcity of patients with ultra-rare disorders, it is infeasible to directly train a model on them. Therefore, we first analyzed the influence of replacing GestaltMatcher DCNN with a state-of-the-art face recognition approach, iResNet with ArcFace. Additionally, we experimented with different face recognition datasets for transfer learning. Furthermore, we proposed test-time augmentation, and model ensembles that mix general face verification models and models specific for verifying disorders to improve the disorder verification accuracy of unseen ultra-rare disorders. Our proposed ensemble model achieves state-of-the-art performance on both seen and unseen disorders. Code is available at github.com/igsb/GestaltMatcher-Arc.

1. Introduction

More than 6% of the global population is affected by rare genetic disorders [11]. Because of the rarity and diversity of genetic disorders, reaching a diagnosis is challenging and time-consuming. More than a third of patients wait for over five years to receive a diagnosis, often referred to as the “diagnostic odyssey” [35]. Many disorders have distinctive

* Corresponding author.

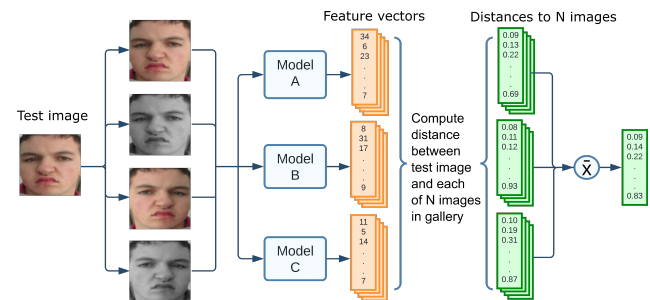


Figure 1. Model ensemble of our approach. We first performed test time augmentation to augment the test image into four images (color and horizontal flip). The four augmented images were further encoded by three different models into 12 representation vectors. We then compared the cosine distance of the 12 representation vectors to the 12 representation vectors from each of the N images in the gallery. It resulted in 12 distance vectors, and each vector contains N cosine distances. In the end, we averaged over 12 distance vectors (\bar{X}) to obtain the final distance vector, which further ranked the N images in the gallery. The gallery image with a smaller distance is more similar to the test image.

dysmorphic facial features, and these features (gestalt) are hints for clinicians to diagnose patients. However, recognizing the facial gestalt presented on a patient’s face highly relies on the clinician’s experience, and it is very difficult if the clinician has never seen the disorder before.

With recent advances in computer vision, many next-generation phenotyping (NGP) approaches have emerged to predict rare disorders by analyzing patient’s 2D frontal image [4, 7, 8, 12, 14, 16, 17, 21, 22, 23, 28, 31, 33]. Among them, DeepGestalt [14] utilized transfer learning to train a deep convolutional neural network on CASIA [34] and to further fine-tune on over 17,106 patient frontal images with 216 disorders. It achieved 91% of top-10 accuracy on a test set of 502 images with 92 different disorders and even out-

performed human experts. Although DeepGestalt demonstrated extraordinary accuracy in predicting these disorders, it can only classify the disorders it has seen during training, and the trained syndromes are only a tiny proportion of all genetic disorders. If the disorders are ultra-rare or novel, we cannot include them in the model training due to a lack of images. These “unseen” syndromes are often cases in the real world (Supplementary Figure S1). Therefore, a way to support unseen syndromes becomes crucial.

To support unseen syndromes, GestaltMatcher was proposed as an extension of DeepGestalt that takes the feature layer before the classification layer in DeepGestalt as the encoder that learned facial dysmorphic features [17]. It encoded the frontal image into a 320-dimensional representation vector. These representation vectors further spanned a feature space. All patients with genetic disorders can be matched or clustered in this space, no longer being limited to the disorders that are trained on (seen) by the networks.

However, both DeepGestalt and GestaltMatcher used the architecture and dataset for transfer learning proposed by Yi *et al.* [34] in 2014. Since then, many larger face recognition datasets [1, 3, 6] and more advanced architectures and loss functions [6, 9, 15, 24, 32] were proposed that achieved higher performance on the face verification task. Therefore, the first aim of this study was to update the architecture by using iResNet [9] and ArcFace [6], and investigating the influence of using different face datasets for transfer learning.

Moreover, a way to train the model that generates better feature representations for unseen ultra-rare disorders has not yet been studied. Hence, the second aim was to investigate different training settings to understand how we can obtain better feature representations for unseen disorders. Our findings showed that fine-tuning on the disorder dataset improved the seen disorder’s accuracy but was not always beneficial for the unseen disorders. Thus, we proposed a model ensemble to integrate face verification and disorder models to improve performance on both seen and unseen syndromes (Figure 1).

In summary, the contributions in this paper are as follows:

- We analyzed the influence of updating the architecture, loss function, and face dataset used for transfer learning.
- We investigated the training settings to generate better feature representations for unseen ultra-rare disorders.
- Every updated individual model outperformed the GestaltMatcher baseline model by [17].
- We proposed a model ensemble to mix general face verification models and models specific for verifying disorders to improve the disorder verification accuracy of unseen ultra-rare disorders.

All experiments were conducted on the GestaltMatcher

Database (GMDB), which is available to medical-related research communities.

2. Related works

2.1. Next-generation phenotyping

Many rare genetic disorders present recognizable facial features, also called “facial gestalt”. For example, patients with Down syndrome have a distinct facial gestalt. Recognizing the facial gestalt shown in a patient’s face is helpful for clinicians in diagnosing the patient. However, it highly relies on the clinician’s experience. When disorders are ultra-rare or novel, a clinician has very likely not seen the disorders before. Therefore, next-generation phenotyping approaches that analyze patients 2D frontal face photo to facilitate the diagnosis become crucial.

In 2014, Ferry *et al.* utilized the shape and appearance representation vectors derived from trained Active Appearance Models and further constructed a feature space they dubbed the “Clinical Face Phenotype Space” (CFPS) using the representation vectors for disorder classification [12]. They trained the model on 1,363 images of eight syndromes and 1,515 images from healthy individuals, and it was the first study that analyzed a relatively large cohort.

With the rapid development of computer vision, many approaches using deep convolutional neural networks (DCNN) have been proposed. Shukla *et al.* [28] trained AlexNet [20] on the entire face and four different facial regions (top right, top left, bottom right, and bottom left) of LFW [19] and concatenated five representation vectors into one 20,480 dimensional vector. In the end, a support vector machine was used to classify six different disorders. Later in 2019, DeepGestalt [14], which utilized transfer learning to train a DCNN on more than 17,106 patient photos with 216 different disorders, showed a high prediction accuracy that outperformed clinical experts. Hong *et al.* [16] also used transfer learning to fine-tune VGG-16 [29] on 228 children with genetic disorders and 228 healthy children. It performed binary classification (with/without a genetic disorder) that could be used for screening.

However, the prevalence of rare disorders is highly imbalanced (Supplementary Figure S1). The number of disorders with enough photos to be included for training the DCNN are a relatively small proportion of all genetic disorders. Syndromes with very few images or novel disorder were not suitable to the classification methods. Therefore, Marbach *et al.* [26] demonstrated matching two unrelated patients with a novel disease by using facial embeddings encoded by FaceNet [27]. In addition, van der Donk *et al.* [31] concatenated the facial embeddings encoded by a normal face recognition model and model trained by disorders. They further performed a clustering analysis to validate the given cohorts with significant facial gestalt. There-

fore, generating facial embeddings that generalize dysmorphic facial features for unseen ultra-rare disorders is essential for rare disorder analysis.

2.2. DeepGestalt

DeepGestalt was proposed by FDNA Inc., which is considered as the current state-of-the-art disorder classification framework [14]. It uses the architecture proposed by Yi *et al.* [34] trained on CASIA [34] to learn general facial features as a base for transfer learning, to later fine-tune the network on 17,106 patient images with 216 different disorders. The architecture, similar to AlexNet, consists of ten convolutional layers, where every two convolutional layers are followed by a pooling layer, and optimizes a Softmax loss function.

Gurovich *et al.* proposed an ensemble method that first cropped the face into multiple regions. The aforementioned architecture was used to train a model for each of the facial regions. In the end, it aggregated the softmax values obtained from each region to perform the diagnosis. It showed 91% of the top-10 accuracy on a test set of 502 images with 92 disorders. In addition to predicting the disorder, it also demonstrated the ability to classify the subtypes of a disorder.

DeepGestalt is used by thousands of clinicians in their daily diagnosis, and is further integrated into the exome sequencing analysis that facilitates the diagnosis on the molecular level [18]. However, as briefly discussed in the previous section, DeepGestalt does not work on ultra-rare or novel disorders unseen during training. Therefore, GestaltMatcher was proposed to overcome this limitation.

2.3. GestaltMatcher

GestaltMatcher [17] is an extension of the DeepGestalt approach. It used the same architecture and face dataset (CASIA) as a base for transfer learning. After training, it used the last 320-dimensional fully-connected layer before the classification layer as the feature layer, and used it as an encoder that encoded each image into a 320-dimensional representation vector. The representation vectors further spanned a CFPS. In the CFPS, patients with rare disorders can be matched to other similar patients. Moreover, clustering analysis can be performed to analyze the similarity among different disorders. GestaltMatcher has been used in several studies to analyze patient similarities [2, 10, 13].

The advantage of GestaltMatcher is that it is no longer limited to the disorders it has seen during training, and it enables researchers to quantify patient-to-patient or syndrome-to-syndrome similarity. However, GestaltMatcher used the same architecture and pre-trained dataset as DeepGestalt that are relatively outdated. Therefore, a study is required to update the architecture and explore methods that improve the performance of unseen ultra-rare

Dataset	# of images	# of individuals
VGG2 [3]	3.31M	9,131
CASIA [34]	0.49M	10,575
MS1MV2 [6]	5.8M	85K
MS1MV3 [6]	5.1M	93K
Glnt360K [1]	17M	360K

Table 1. Overview of the face datasets.

disorders verification.

3. Datasets and methodology

3.1. Datasets

3.1.1 Face recognition datasets

In this paper, we experimented with five different face recognition datasets to be used for training the (transfer learning) base model: VGG2 [3], CASIA [34], MS1MV2 [6], MS1MV3 [6], and Glnt360K [1]. The full name of CASIA dataset is CASIA-WebFace. We used CASIA as abbreviation in this paper. The number of images in the datasets ranges from 0.49M to 17M. An overview of the datasets is shown in Table 1.

3.1.2 GestaltMatcher Database - rare disorder dataset

Hsieh *et al.* [17] built up GestaltMatcher Database¹ (GMDB), which collects medical images of rare disorder from publications and patients with proper consent from clinics. It is open to clinicians and researchers working in medical research fields. To avoid data abuse, applicants need to be reviewed by a committee of GMDB before they can access the database.

We used GMDB (v1.0.3) to fine-tune the base models on faces of patients with disorders. GMDB (v1.0.3) contains 7,459 frontal images of 5,995 patients with 449 different disorders. All the disorders have at least two patients. The dataset was further divided into two sets, a “frequent” (GMDB-Frequent) and a “rare” (GMDB-Rare) set. The disorders with more than six patients were assigned to GMDB-Frequent, while the disorders with six or fewer patients were assigned to GMDB-Rare.

There are 6,354 images of 5,123 patients with 204 disorders in GMDB-Frequent and 1,105 images of 872 patients with 245 disorders in GMDB-Rare. We fine-tuned the base models on GMDB-Frequent, thus disorders in this set can be considered as “seen” disorders. For training, GMDB-Frequent was further divided into 5,100 images for the training set, 661 images for the validation set, and 593 images for the test set. On the other hand, GMDB-Rare was “unseen” during training. We used GMDB-Rare to simulate

¹<https://db.gestaltmatcher.org/>

Dataset	# of images	# of patients	# of disorders
GMDB-Frequent	6,354	5,123	204
GMDB-Rare	1,105	872	245
Total	7,459	5,995	449

Table 2. Overview of GMDB dataset. GMDB-Frequent is used for fine-tuning and thus “seen” by the model, while disorders in GMDB-Rare are “unseen” to the model.

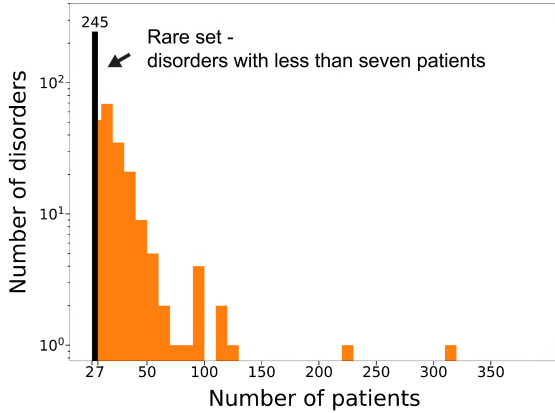


Figure 2. Disorder distribution of GMDB. The X-axis shows the number of patients in the disorder. The Y-axis shows the number of disorders with the corresponding number of patients in X-axis, and it is in log scale. The black bar is the rare set (GMDB-Rare), that has disorders with more than one patient and fewer than seven patients.

ultra-rare or novel disorders in real-world scenarios. The overview of the GMDB dataset is shown in Table 2. In Figure 2, GMDB shows a long tail distribution. GMDB-Rare has only 14.5% (872/5995) of all patients, but it covers 54.5% of the disorders. The distribution of GMDB is similar to the estimation of disorder prevalence in the real world (Supplementary Figure S1).

3.2. Evaluation

We evaluated the performance of the base models on the popular face verification dataset Labeled Faces in the Wild (LFW) [19]. During evaluation, two faces were compared, and the models verified whether they belong to the same person. The evaluation set consisted of 11-folds, where the first was used to establish a threshold, while the remaining 10-folds were used for the final evaluation.

Most importantly, we used GMDB to evaluate the base models, our fine-tuned models, and eventually the model ensemble. During the evaluation procedure, the feature space was populated with a gallery set of diagnosed patients’ feature vectors, being either the seen disorders

(GMDB-Frequent), the unseen disorders (GMDB-Rare) or an unified gallery (GMDB-Frequent and -Rare). Afterwards, representation vectors of test images were matched to the gallery cases in the feature space. For the unseen test images, 10-fold cross-validation was performed.

We first calculated the cosine distances between the test image and each image in the gallery. The cosine distance further ranked the gallery images. Hsieh *et al.* [17] showed the top- k ($k \in [1, 5, 10, 30]$) mean accuracy (as described in Equation 1) for test images of seen and unseen disorders. Instead, we focused on the top-1 and top-5 results in this paper, though top-10 and top-30 are included in the Supplemental Tables. The GestaltMatcher DCNN² (GM-Hsieh2022) was used as the baseline, having retrained it on the recent version of GMDB.

$$mA_k = \frac{1}{C} \sum_c A_{k,c}, \quad (1)$$

where mA_k is the top- k mean accuracy, C is the number of classes, c is the class index, and $A_{k,c}$ is the top- k accuracy for class c .

3.3. Model architecture and training

Our model architecture is based on the one used by Deng *et al.* [6]. They used a popular variation on the ResNet architecture named iResNet [9]. It includes more batch normalization. In addition to the original implementation, it replaces the ReLU activation function with PReLU, and lastly, it batch normalizes the computed representation vector.

The training procedure was split into two steps:

- Training a base model on a face recognition dataset for transfer learning;
- Fine-tuning that base model on GMDB for disorder verification.

For the first part, we used pre-trained models supplied by insightface³ that have been trained on different face recognition datasets using Additive Angular Margin Loss (ArcFace). The loss is defined by the Equation 2.

$$L_{Arc} = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{s \cos(\theta_{y_i} + m)}}{e^{s \cos(\theta_{y_i} + m)} + \sum_{j=1, j \neq y_i}^n e^{s \cos \theta_j}}, \quad (2)$$

where θ_j is the angle between weights W_j and feature x_i . We insert angular margin m to get a new angle between the true logit, y_i , and our representation vector to become $\theta_{y_i} + m$. s is the scale of L_2 normalized representation vectors. The pre-trained models used m and s set to 0.5 and 64, respectively. Representations learned using this loss tended

²<https://github.com/igsb/GestaltMatcher>

³<https://github.com/deepinsight/insightface>

to have stronger distinctions between different classes and better similarities for the same classes than most other metric learning losses.

An important preprocessing step of insightface’s training procedure is to align faces based on five landmarks: left and right eye, nose, left and right mouth corner. This alignment is essential to reproduce their performance. For our implementation, we used RetinaFace [5] to obtain these landmarks and the alignment code they supplied, which uses an affine transformation based on matching landmark locations.

For the second part, we made minor changes to the model architecture. We removed the batch normalization of the computed features. This normalization was necessary for ArcFace, which we did not use during fine-tuning. Instead, due to the small dataset size, significant class imbalance, and long tail distributions, we decided only to optimize Weighted Cross Entropy Softmax Loss (WCE).

To address the class imbalance, we used Equation 3 to calculate the class weights, casting them into the range $(0.5, \dots, 1.0]$.

$$W_c = \frac{0.5 \cdot \min(D)}{D_c} + 0.5, \quad (3)$$

where D is the set of frequencies per class, c is the class, and W_c is the WCE weight for class c . If not for our lower bound of $W_c > 0.5$, due to the long-tailed distribution it would be possible for W_c to be lower than 0.01. This would make the training process challenging. We also replaced the final fully connected layer to train a classifier on the disorders of the training set. Lastly, we froze all model weights except those of the feature and classification layer.

We fine-tuned our model on GMDB using aligned faces of size 112x112, randomly flipping horizontally, randomly converting color images to gray, color jittering, and randomly adding zooming/cropping artifacts. We further used the Adam optimizer with a base learning rate of 1e-3, which was reduced by a factor of 2 when the top-5 mean accuracy on the validation set plateaus until convergence. The mean accuracy was calculated with Equation 1. Code is available at github.com/igsb/GestaltMatcher-Arc.

3.4. Inference strategy

An essential part of our method was our strategy during inference. We aimed to improve our performance on both seen and unseen disorders by computing multiple representation vectors per image, aiming to end up with a better overall ranking than each separate representation vector. We employed two approaches to obtain multiple representation vectors per image: model ensembles and test time augmentation.

3.4.1 Model ensembles

Model ensembles are mixtures of models that combine each model’s output. This approach helps achieve a better overall generalization as it leverages each model’s strengths to alleviate the others’ weaknesses. In our case, we presented each model with the same image, computed each model’s representation vector, and averaged the cosine distances from the image to the GMDB gallery set.

For our ensemble, we considered both models that are fine-tuned for disorders and models built for face verification. The face verification models produced strong general features that can be leveraged to verify unseen disorders, while the fine-tuned models were fitted towards features of seen disorders they have been trained on. More specifically, we included one face verification model, one deeper disorder model (using iResNet-100), and one disorder model (using iResNet-50) designed to be less prone to overfitting on the seen disorders. More detailed information on the selected models can be found in Section 4.

3.4.2 Test time augmentation

Test time augmentation (TTA), similarly to model ensembles, combines outputs to achieve more robust performance. However, instead of presenting different models with the same image, it presented the same model with an image and augmented versions of that image (e.g., horizontally flipped, converted from color to gray, rotation, and translation). Ideally, the representation vectors would be close to identical because they are from the same image, and the actual face does not change. In practice, this is usually not the case. This helps average out the cosine distance between the gallery and test set.

Of course, not all augmentations make sense to use during TTA. Any augmentation that changes the face structure or affects the required face alignment is potentially harmful for our implementation. Generally, augmentations used during training are well suited for TTA. As such, we used horizontal flipping and conversion from color to gray.

Finally, we averaged the cosine distance of all models in the ensemble and TTAs per model (i.e., three models and two TTA each, $3 \times 2 \times 2 = 12$ cosine distances). The order of the disorders ranked for verification was determined by the k -nearest neighbors of the average cosine distance between the gallery images to the test image. We decided to use $k = 1$ due to the highly imbalanced data in GMDB (and thus also in the gallery set). Using a $k > 1$ would be problematic for disorders with only one occurrence in the gallery set. Figure 1 gives a simplistic view of how these inference strategies work.

4. Experiments and results

4.1. Updating the architecture and base optimization

We hypothesized that replacing GestaltMatcher’s base model (GM-Hsieh2022) with a state-of-the-art face verification model will improve the overall performance on LFW and GMDB.

First, we compared the performance of the GM-Hsieh2022 model, using an AlexNet-like architecture and cross-entropy loss, to iResNet-50 with ArcFace, both trained on CASIA. Afterward, we fine-tuned these models on GMDB. During the fine-tuning process, both models only used weighted cross entropy. Results of the base models and fine-tuned models are shown in Table 3. An extended version of the table can be found in Supplementary Table S1, and the performance when using the unified gallery in Supplementary Table S5.

We find that the features generated by the ArcFace base model are generally more descriptive than those of the GM-Hsieh2022 base model. This is supported by the higher LFW performance and the overall higher performance on GMDB without fine-tuning. In Table 3, the LFW accuracy is increased from 93.8% to 98.4%, and the top-1 and top-5 accuracies for both seen (GMDB-Frequent) and unseen (GMDB-Rare) disorders are improved when we update the model from GM-Hsieh2022 to ArcFace-r50.

After fine-tuning, the GM-Hsieh2022 model improved on both GMDB-Frequent and GMDB-Rare. The ArcFace model significantly increased the performance on seen disorders while decreasing the performance on unseen disorders. We believed this indicated that the model has a higher tendency to overfit on the small dataset than the GM-Hsieh2022 model had. Although the performance of GMDB-Rare dropped after fine-tuning the new model (ArcFace-r50*), the top-1 and top-5 accuracies were still similar to the fine-tuned GestaltMatcher by [17] (GM-Hsieh2022*).

Model	LFW	GMDB-Frequent		GMDB-Rare	
		Top-1	Top-5	Top-1	Top-5
GM-Hsieh2022	93.8%	10.99%	29.39%	14.64%	27.03%
GM-Hsieh2022*	-	15.96%	33.83%	19.26%	36.28%
ArcFace-r50	98.4%	21.84%	40.87%	22.74%	37.35%
ArcFace-r50*	-	35.37%	53.25%	19.29%	36.00%

Table 3. Comparison of the performance of the GM-Hsieh2022 model and the ArcFace-r50 model on LFW and GMDB. Both have been pre-trained on CASIA and models marked with (*) have been fine-tuned on GMDB. For each column, the best accuracy between the models before fine-tuning and after fine-tuning is boldfaced.

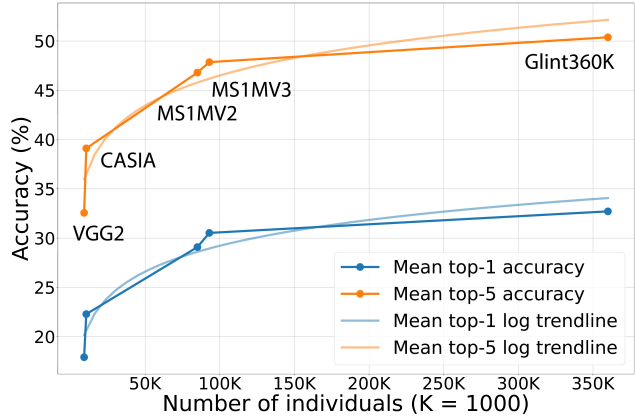


Figure 3. Mean accuracy of the ArcFace-r50 base model on GMDB when using different datasets. The X-axis shows the number of individuals in the datasets. The Y-axis shows the mean accuracy (both GMDB-Frequent and -Rare) of models using different base datasets. The light orange and blue lines shows the logarithmic relation.

4.2. Updating the transfer learning dataset

We hypothesized that increasing the number of individuals in the transfer learning base dataset will result in better/more general (facial) feature descriptors. However, we expected some drop-off in the performance gain when increasing the number of individuals indefinitely. To test this hypothesis, we compared the performance on LFW and GMDB with five well-known face recognition datasets: VGG2, CASIA, MS1MV2, MS1MV3, and Glint360K. The results are shown in Table 4 and Figure 3. An extended version of the table can be found in Supplementary Table S2, and the performance when using the unified gallery in Supplementary Table S6.

Figure 3 shows the average accuracy per base dataset concerning the number of unique individuals in the dataset. Table 4 shows that the ArcFace-r50 base models trained on datasets with more different individuals tends to achieve higher accuracy on both GMDB-Frequent and GMDB-Rare before fine-tuning. We also found a drop-off in accuracy gained when increasing the number of individuals based on the logarithmic relation shown in the Figure 3.

Moreover, we found that the accuracy on GMDB-Frequent after fine-tuning did not always improve when using a larger dataset. For example, in Table 4, the top-1 and top-5 accuracies of Glint360K* are lower than the accuracies of MS1MV3*, which drop from 45.06% to 41.58% and 64.64% to 62.60%, respectively. However, the accuracy on GMDB-Rare after fine-tuning always improved when we used a larger dataset for training the ArcFace-r50 base model.

The results show that both GMDB-Frequent and

Dataset	LFW	GMDB-Frequent		GMDB-Rare	
		Top-1	Top-5	Top-1	Top-5
VGG2	98.5%	15.52%	31.56%	20.31%	33.57%
CASIA	98.4%	21.84%	40.87%	22.74%	37.35%
MS1MV2	99.0%	29.14%	48.86%	29.04%	44.74%
MS1MV3	98.9%	31.54%	49.36%	29.52%	46.36%
Glint360K	99.0%	32.43%	53.14%	33.00%	47.62%
VGG2*	85.8% (-12.7%)	27.50% (+11.98%)	49.92% (+18.36%)	17.56% (-2.75%)	33.41% (-0.16%)
CASIA*	75.7% (-22.7%)	35.37% (+13.53%)	53.25% (+12.38%)	19.29% (-3.45%)	36.00% (-1.35%)
MS1MV2*	84.1% (-17.7%)	39.98% (+10.84%)	59.81% (+10.95%)	21.86% (-7.18%)	39.89% (-4.85%)
MS1MV3*	76.4% (-22.5%)	45.06% (+13.52%)	64.64% (+15.28%)	24.31% (-5.21%)	40.28% (-6.08%)
Glint360K*	84.9% (-12.6%)	41.58% (+9.15%)	62.60% (+9.46%)	26.55% (-6.45%)	42.69% (-4.93%)

Table 4. Comparison of the performance of the ArcFace-r50 models trained on a variety of face recognition datasets. The percentages within parentheses indicate the change between the face verification (base) and the fine-tuned models. For example, the top-1 accuracy on GMDB-Frequent increased by 11.98% (15.52% \rightarrow 27.50%) from VGG2 to VGG2*. Models marked with (*) have been fine-tuned on GMDB.

Model	LFW	GMDB-Frequent		GMDB-Rare	
		Top-1	Top-5	Top-1	Top-5
r50	84.9%	41.58%	62.60%	26.55%	42.69%
r50-D/O	86.2%	46.95%	66.07%	28.85%	45.36%
r50-D/O†	87.6%	44.33%	65.76%	29.06%	46.35%
r100	91.0%	47.96%	68.87%	26.03%	42.22%
r100-D/O	91.1%	48.37%	71.78%	28.02%	44.32%
r100-D/O†	93.0%	49.25%	69.95%	30.33%	47.85%

Table 5. Comparison of the performance of iResNet-50 and -100 fine-tuned on GMDB. D/O indicates an additional dropout layer and (†) indicates the use of L_2 weight decay on the feature layer. For each column, the best accuracy among the models (without regularization, D/O, and D/O†) is boldfaced.

GMDB-Rare benefit from using a larger dataset for training the ArcFace-r50 base model, especially for the unseen disorders (GMDB-Rare). In addition, it might not be necessary to use a face recognition dataset larger than Glint360K, as the performance gain seems to saturate when the number of individuals in the dataset is larger than 1M.

4.3. Influence of fine-tuning ArcFace on GMDB

In an earlier experiment, we saw that fine-tuning ArcFace on GMDB reduced the accuracy on unseen disorders (GMDB-Rare). We believed that fine-tuning the feature representation layer on GMDB will negatively influence the general feature descriptors’ quality by (over)fitting on the small imbalanced dataset, not just when using CASIA as the base dataset but also for larger base datasets. We believed this should reflect in the accuracy on LFW. As such, we fine-tuned the ArcFace models trained on the CASIA, VGG2, MS1MV2, MS1MV3, and Glint360K on GMDB, and afterward evaluated them on LFW and GMDB. The results are shown in Table 4. An extended version of the

table can be found in Supplementary Table S2, and the performance when using the unified gallery in Supplementary Table S6.

Based on the results in the Table 4, we find that fine-tuning decreases the performance on unseen disorders (GMDB-Rare) for every model, as well as the general face verification performance on LFW. However, the performance of the models using a larger base dataset still outperform the baseline for these unseen disorders.

4.4. Additional regularization during fine-tuning to improve generalizability on unseen disorders

We believed that the overfitting, shown in the previous experiment as a drop in accuracy for unseen disorders, can be reduced by adding additional regularization to the feature layer in the form of L_2 weight decay and dropout. We fine-tuned the iResNet-50 and iResNet-100 ArcFace models pre-trained on Glint360K to include additional dropout and additional L_2 weight decay of the feature layer ($\lambda = 5e^{-5}$). The results are shown in Table 5. An extended version of the table can be found in Supplementary Table S3, and the performance when using the unified gallery in Supplementary Table S7.

We find that the use of additional dropout improves accuracy for both the seen and unseen disorders. Additional L_2 weight decay on the feature layer helps to maintain some of the LFW performance and, on some occasions, is an improvement to dropout. For example, the top-1 accuracy on GMDB-Frequent increases from 48.37% to 49.25% when applying L_2 weight decay to r100-D/O. Although the improvement of L_2 weight decay on seen disorders (GMDB-Frequent) is inconclusive, the improvement on the unseen ones (GMDB-Rare) is clear.

Model	Dataset	Loss	GMDB-Frequent		GMDB-Rare	
			Top-1	Top-5	Top-1	Top-5
GM-Hsieh2022	CASIA*	CE	15.96%	33.83%	19.26%	36.28%
r50-D/O†	Glnt360K*	CE	44.33%	65.76%	29.06%	46.35%
r50-D/O†+ TTA	Glnt360K*	CE	47.73%	67.67%	30.29%	46.38%
r100-D/O	Glnt360K*	CE	48.37%	71.78%	28.02%	44.32%
r100-D/O + TTA	Glnt360K*	CE	51.16%	69.58%	27.92%	46.26%
r100	Glnt360K	ArcFace	30.25%	54.81%	33.25%	50.22%
r100 + TTA	Glnt360K	ArcFace	35.25%	56.52%	33.47%	51.61%
Model ensemble	n/a	n/a	52.06%	70.70%	34.93%	52.78%
Model ensemble + TTA	n/a	n/a	52.99%	71.01%	35.98%	53.93%

Table 6. Comparison of the performance of the GM-Hsieh2022 (baseline) model, two ArcFace models fine-tuned on GMDB, one ArcFace face verification model, and our model ensemble using the three ArcFace models. TTA indicates the model was evaluated using test time augmentation, (*) indicates the model was fine-tuned on GMDB, (D/O) indicates an additional dropout layer, and (†) indicates the use of L_2 weight decay on the feature layer.

4.5. Influence of inference strategies

We believed the inference strategies discussed in Section 3.4 will improve most models’ accuracy. We hypothesized that presenting our model with an image and slight variations of that image will increase the robustness of the clustering. On top of that, combining our disorder models, fine-tuned on GMDB, with general face verification models will improve generalizability and robustness for both seen and unseen disorders. Table 6 shows the performance of the baseline (GM-Hsieh2022), each model used in the ensemble with and without TTA, and the model ensemble with and without TTA. An extended version of the table can be found in Supplementary Table S4, and the performance when using the unified gallery in Supplementary Table S8.

In Table 6, we find that TTA increases almost every test group’s performance. Only the top-5 accuracy on GMDB-Frequent and top-1 accuracy on GMDB-Rare for r100 with dropout (r100-D/O) are decreased after applying TTA. Moreover, the model ensemble outperforms every single model in almost every test group, except the top-5 accuracy on GMDB-Frequent. The top-5 accuracy on GMDB-Frequent for the model ensemble is 70.70% which is slightly lower than 71.78% from r100 with dropout (r100-D/O). In the end, combining the model ensemble and TTA further improves the performance, achieving state-of-the-art. When comparing the model ensemble with TTA to the GM-Hsieh2022 model, the top-1 accuracy improves from 15.96% to 52.99% and 19.26% to 35.98% on GMDB-Frequent and GMDB-Rare, respectively, showing a strong performance on both seen and unseen disorders.

5. Conclusion and future works

We found that using face recognition datasets with more individuals led to more generalized representation vectors,

which in turn form a good base for transfer learning. Fine-tuning the transfer learning datasets with ArcFace iResNet on GMDB led to a significant increase in performance on seen disorders and a decrease in performance on unseen disorders. The latter is likely caused by overfitting on the seen disorders, which unlearns general facial features. The use of regularization techniques such as dropout and L_2 weight decay can help reduce the impact of overfitting, increasing the performance of unseen disorders. Moreover, using TTA increased the performance of all models. Next, combining one face verification model and two disorder verification models in a model ensemble allowed us to leverage their strengths on both seen and unseen disorders.

In conclusion, each model with and without TTA, and our model ensemble outperformed GM-Hsieh2022, where the model ensemble achieves state-of-the-art performance. We believe this work can function as a strong baseline for future comparison in this emerging field.

In this study, we focused on the imbalanced number of patients among the disorders. However, Lumaka *et al.* reported that the performance of DeepGestalt was biased by the imbalance of ethnicity groups in the training set [25]. Therefore, an approach that consider the imbalance of ethnicity, sex, and age is important.

Moreover, we only discussed the iResNet with ArcFace and cross entropy. Benchmarking on different architectures and loss functions, such as EfficientNet [30], CosFace [32], and SphereFace [24] is required to understand more on how to obtain more generalized representation vectors for unseen disorders. Besides, using different representation vector dimensions and dimension reduction methods are also a possibility to further optimize the feature representation for unseen disorders.

References

- [1] Xiang An, Xuhan Zhu, Yuan Gao, Yang Xiao, Yongle Zhao, Ziyong Feng, Lan Wu, Bin Qin, Ming Zhang, Debing Zhang, and Ying Fu. Partial FC: Training 10 million identities on a single machine. In *2021 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, pages 1445–1449, Oct. 2021.
- [2] Maria Asif, Emrah Kaygusuz, Marwan Shinawi, Anna Nickelsen, Tzung-Chien Hsieh, Prerana Wagle, Birgit Budde, Jennifer Hochscherf, Uzma Abdullah, Stefan Höning, Christian Nienberg, Dirk Lindenblatt, Angelika A Noegel, Janine Altmüller, Holger Thiele, Susanne Motameny, Nicole Fleischer, Idan Segal, Lynn Pais, Sigrid Tinschert, Nadra G Samra, Juliann M Savatt, Natasha L Rudy, Chiara De Luca, Paola Fortugno, Susan M White, Peter Krawitz, Anna C E Hurst, Karsten Niefind, Joachim Jose, Francesco Brancati, Peter Nürnberg, and Muhammad Sajid Hussain. De novo variants of CSNK2B cause a new intellectual disability-craniodigital syndrome by disrupting the canonical wnt signaling pathway. *Human Genetics and Genomics Advances*, page 100111, Apr. 2022.
- [3] Qiong Cao, Li Shen, Weidi Xie, Omkar M Parkhi, and Andrew Zisserman. VGGFace2: A dataset for recognising faces across pose and age. Oct. 2017.
- [4] J J Cerrolaza, A R Porras, A Mansoor, Q Zhao, M Summar, and M G Linguraru. Identification of dysmorphic syndromes using landmark-specific local texture descriptors. In *2016 IEEE 13th International Symposium on Biomedical Imaging (ISBI)*, pages 1080–1083, Apr. 2016.
- [5] Jiankang Deng, Jia Guo, Evangelos Ververas, Irene Kotzia, and Stefanos Zafeiriou. RetinaFace: Single-Shot Multi-Level face localisation in the wild. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5202–5211, June 2020.
- [6] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. ArcFace: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 2019-June, pages 4685–4694. IEEE Computer Society, June 2019.
- [7] Tracy Dudding-Byth, Anne Baxter, Elizabeth G Holliday, Anna Hackett, Sheridan O’Donnell, Susan M White, John Attia, Han Brunner, Bert de Vries, David Koolen, Tjitske Kleefstra, Seshika Ratwate, Carlos Riveros, Steve Brain, and Brian C Lovell. Computer face-matching technology using two-dimensional photographs accurately matches the facial gestalt of unrelated individuals with the same syndromic form of intellectual disability. *BMC Biotechnol.*, 17(1):1–9, 2017.
- [8] Dat Duong, Ping Hu, Cedrik Tekendo-Ngongang, Suzanna E Ledgister Hanchard, Simon Liu, Benjamin D Solomon, and Rebekah L Waikel. Neural networks for classification and image generation of aging in genetic syndromes. *Front. Genet.*, 13, 2022.
- [9] Ionut Cosmin Duta, Li Liu, Fan Zhu, and Ling Shao. Improved residual networks for image and video recognition. pages 9415–9422, Jan. 2021.
- [10] Frédéric Ebstein, Sébastien Küry, Victoria Most, Cory Rosenfelt, Marie-Pier Scott Boyer, Geeske M van Woerden, Thomas Besnard, Jonas Johannes Papendorf, Maja Studencka-Turski, Tianyun Wang, Tzung-Chien Hsieh, Richard Golnik, Dustin Baldrige, Cara Forster, Charlotte de Konink, Selina M W Teurlings, Virginie Vignard, Richard H van Jaarsveld, Lesley Ades, Benjamin Cogné, Cyril Mignot, Wallid Deb, Marjolijn C J Jongmans, F Sessions Cole, Marie-José H van den Boogaard, Jennifer A Wambach, Daniel J Wegner, Sandra Yang, Vickie Hannig, Jennifer Ann Brault, Neda Zadeh, Bruce Bennetts, Boris Keren, Anne-Claire Gelineau, Zöe Powis, Meghan Towne, Kristine Bachman, Andrea Seeley, Anita E Beck, Jennifer Morrison, Rachel Westman, Kelly Averill, Theresa Brunet, Judith Haasters, Melissa T Carter, Matthew Osmond, Patricia G Wheeler, Francesca Forzano, Shehla Mohammed, Yannis Trakadis, Andrea Accogli, Rachel Harrison, Sophie Rondeau, Geneviève Baujat, Giulia Barcia, René Günther Feichtinger, Johannes Adalbert Mayr, Martin Preisel, Frédéric Laumonier, Alexej Knaus, Bertrand Isidor, Peter Krawitz, Uwe Völker, Elke Hammer, Arnaud Droit, Evan E Eichler, Ype Elgersma, Peter W Hildebrand, François Bolduc, Elke Krüger, Stéphane Bézieau, Deciphering Developmental Disorders Study, and Care4Rare Canada Consortium. De novo variants in the PSMC3 proteasome AAA-ATPase subunit gene cause neurodevelopmental disorders associated with type I interferonopathies. Dec. 2021.
- [11] Carlos R Ferreira. The burden of rare diseases. *Am. J. Med. Genet. A*, 179(6):885–892, June 2019.
- [12] Quentin Ferry, Julia Steinberg, Caleb Webber, David R Fitzpatrick, Chris P Ponting, Andrew Zisserman, and Christoffer Nellåker. Diagnostically relevant facial gestalt information from ordinary photos. *Elife*, 3:e02020, June 2014.
- [13] Lily Guo, Jiyeon Park, Edward Yi, Elaine Marchi, Tzung-Chien Hsieh, Yana Kibalnyk, Yolanda Moreno-Sáez, Saskia Biskup, Oliver Puk, Carmela Beger, Quan Li, Kai Wang, Anastassia Voronova, Peter M Krawitz, and Gholson J Lyon. KBG syndrome: videoconferencing and use of artificial intelligence driven facial phenotyping in 25 new patients. *Eur. J. Hum. Genet.*, pages 1–11, Aug. 2022.
- [14] Yaron Gurovich, Yair Hanani, Omri Bar, Guy Nadav, Nicole Fleischer, Dekel Gelbman, Lina Basel-Salmon, Peter M Krawitz, Susanne B Kamphausen, Martin Zenker, Lynne M Bird, and Karen W Gripp. Identifying facial phenotypes of genetic disorders using deep learning. *Nat. Med.*, 25(1):60–64, Jan. 2019.
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [16] Dian Hong, Ying-Yi Zheng, Ying Xin, Ling Sun, Hang Yang, Min-Yin Lin, Cong Liu, Bo-Ning Li, Zhi-Wei Zhang, Jian Zhuang, Ming-Yang Qian, and Shu-Shui Wang. Genetic syndromes screening by facial recognition technology: VGG-16 screening model construction and evaluation. *Orphanet J. Rare Dis.*, 16(1):344, Aug. 2021.
- [17] Tzung-Chien Hsieh, Aviram Bar-Haim, Shahida Moosa, Nadja Ehmke, Karen W Gripp, Jean Tori Pantel, Mag-

- dalena Danyel, Martin Atta Mensah, Denise Horn, Stanislav Rosnev, Nicole Fleischer, Guilherme Bonini, Alexander Hustinx, Alexander Schmid, Alexej Knaus, Behnam Javanmardi, Hannah Klinkhammer, Hellen Lesmann, Sugirthan Sivalingam, Tom Kamphans, Wolfgang Meiswinkel, Frédéric Ebstein, Elke Krüger, Sébastien Küry, Stéphane Bézieau, Axel Schmidt, Sophia Peters, Hartmut Engels, Elisabeth Mangold, Martina Kreiß, Kirsten Cremer, Claudia Perne, Regina C Betz, Tim Bender, Kathrin Grundmann-Hauser, Tobias B Haack, Matias Wagner, Theresa Brunet, Heidi Beate Bentzen, Luisa Averdunk, Kimberly Christine Coetzer, Gholson J Lyon, Malte Spielmann, Christian P Schaaf, Stefan Mundlos, Markus M Nöthen, and Peter M Krawitz. GestaltMatcher facilitates rare disease matching using facial phenotype descriptors. *Nat. Genet.*, Feb. 2022.
- [18] Tzung Chien Hsieh, Martin A Mensah, Jean T Pantel, Dione Aguilar, Omri Bar, Allan Bayat, Luis Becerra-Solano, Heidi B Bentzen, Saskia Biskup, Oleg Borisov, Oivind Braaten, Claudia Ciaccio, Marie Coutelier, Kirsten Cremer, Magdalena Danyel, Svenja Daschkey, Hilda David Eden, Koenraad Devriendt, Sandra Wilson, Sofia Douzgou, Dejan ukić, Nadja Ehmke, Christine Fauth, Björn Fischer-Zirnsak, Nicole Fleischer, Heinz Gabriel, Luitgard Graul-Neumann, Karen W Gripp, Yaron Gurovich, Asya Gusina, Nechama Haddad, Nurulhuda Hajjir, Yair Hanani, Jakob Hertzberg, Konstanze Hoernagel, Janelle Howell, Ivan Ivanovski, Angela Kaindl, Tom Kamphans, Susanne Kamphausen, Catherine Karimov, Hadil Kathom, Anna Keryan, Alexej Knaus, Sebastian Köhler, Uwe Kornak, Alexander Lavrov, Maximilian Leitheiser, Gholson J Lyon, Elisabeth Mangold, Purificación Marín Reina, Antonio Martínez Carrascal, Diana Mitter, Laura Morlan Herrador, Guy Nadav, Markus Nöthen, Alfredo Orrico, Claus Eric Ott, Kristen Park, Borut Peterlin, Laura Pölsler, Annick Raas-Rothschild, Linda Randolph, Nicole Revencu, Christina Ringmann Fagerberg, Peter Nick Robinson, Stanislav Rosnev, Sabine Rudnik, Gorazd Rudolf, Ulrich Schatz, Anna Schossig, Max Schubach, Or Shanon, Eamonn Sheridan, Pola Smirin-Yosef, Malte Spielmann, Eun Kyung Suk, Yves Sznajer, Christian T Thiel, Gundula Thiel, Alain Verloes, Irena Vrekar, Dagmar Wahl, Ingrid Weber, Korina Winter, Marzena Wiśniewska, Bernd Wollnik, Ming W Yeung, Max Zhao, Na Zhu, Johannes Zschocke, Stefan Mundlos, Denise Horn, and Peter M Krawitz. PEDIA: prioritization of exome data by image analysis. *Genet. Med.*, 21(12):2807–2814, Dec. 2019.
- [19] Gary B Huang, Manu Ramesh, Tamara Berg, and Erik Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical report, University of Massachusetts, Amherst, Oct. 2007.
- [20] Krizhevsky, Sutskever, and others. Imagenet classification with deep convolutional neural networks. *Adv. Neural Inf. Process. Syst.*
- [21] Kaya Kuru, Mahesan Niranjan, Yusuf Tunca, Erhan Osvank, and Tayyaba Azim. Biomedical visual data analysis to build an intelligent diagnostic decision support system in medical genetics. *Artif. Intell. Med.*, 62(2):105–118, Oct. 2014.
- [22] T Liehr, N Acquarola, K Pyle, S St-Pierre, M Rinholm, O Bar, K Wilhelm, and I Schreyer. Next generation phenotyping in emanuel and Pallister-Killian syndrome using computer-aided facial dysmorphology analysis of 2D photos. *Clin. Genet.*, 93(2):378–381, 2018.
- [23] Hui Liu, Zi-Hua Mo, Hang Yang, Zheng-Fu Zhang, Dian Hong, Long Wen, Min-Yin Lin, Ying-Yi Zheng, Zhi-Wei Zhang, Xiao-Wei Xu, Jian Zhuang, and Shu-Shui Wang. Automatic facial recognition of Williams-Beuren syndrome based on deep convolutional neural networks. *Front Pediatr.*, 9:648255, May 2021.
- [24] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. SphereFace: Deep hypersphere embedding for face recognition. *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, 2017-Janua:6738–6746, 2017.
- [25] A Lumaka, N Cosemans, A Lulebo Mampasi, G Mubungu, N Mvuama, T Lubala, S Mbuyi-Musanazayi, J Breckpot, M Holvoet, T de Ravel, G Van Buggenhout, H Peeters, D Donnai, L Mutesa, A Verloes, P Lukusa Tshilobo, and K Devriendt. Facial dysmorphism is influenced by ethnic background of the patient and of the evaluator. *Clin. Genet.*, 92(2):166–171, Aug. 2017.
- [26] Felix Marbach, Cecilie F Rustad, Angelika Riess, Dejan ukić, Tzung Chien Hsieh, Itamar Jobani, Trine Prescott, Andrea Bevot, Florian Erger, Gunnar Houge, Maria Redfors, Janine Altmueller, Tomasz Stokowy, Christian Gilissen, Christian Kubisch, Emanuela Scarano, Laura Mazzanti, Torunn Fiskerstrand, Peter M Krawitz, Davor Lessel, and Christian Netzer. The discovery of a LEMD2-Associated nuclear envelopathy with early progeroid appearance suggests advanced applications for AI-Driven facial phenotyping. *Am. J. Hum. Genet.*, 104(4):749–757, Apr. 2019.
- [27] Florian Schroff, Dmitry Kalenichenko, and James Philbin. FaceNet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 07-12-June-2015, pages 815–823. IEEE Computer Society, Oct. 2015.
- [28] P Shukla, T Gupta, A Saini, P Singh, and R Balasubramanian. A deep learning Frame-Work for recognizing developmental disorders. In *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 705–714, Mar. 2017.
- [29] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for Large-Scale image recognition. Sept. 2014.
- [30] Tan and Le. Efficientnet: Rethinking model scaling for convolutional neural networks. *International conference on machine learning*.
- [31] Roos van der Donk, Sandra Jansen, Janneke H M Schuurs-Hoeijmakers, David A Koolen, Lia C M J Goltstein, Alexander Hoischen, Han G Brunner, Patrick Kemmeren, Christoffer Nellåker, Lisenka E L M Vissers, Bert B A de Vries, and Jayne Y Hehir-Kwa. Next-generation phenotyping using computer vision algorithms in rare genomic neurodevelopmental disorders. *Genet. Med.*, 21(8):1719–1725, Aug. 2019.

- [32] Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Dihong Gong, Jingchao Zhou, Zhifeng Li, and Wei Liu. CosFace: Large margin cosine loss for deep face recognition. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 5265–5274, 2018.
- [33] Kuan Wang and Jiebo Luo. Detecting visually observable disease symptoms from faces. *EURASIP J. Bioinform. Syst. Biol.*, 2016(1):13, Dec. 2016.
- [34] Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z Li. Learning face representation from scratch. Nov. 2014.
- [35] Yvonne Zurynski, Marie Deverell, Troy Dalkeith, Sandra Johnson, John Christodoulou, Helen Leonard, Elizabeth J Elliott, and APSU Rare Diseases Impacts on Families Study group. Australian children living with rare diseases: experiences of diagnosis and perceived consequences of diagnostic delays. *Orphanet J. Rare Dis.*, 12(1):68, Apr. 2017.