

1 **GestaltMatch: breaking the limits of rare disease matching**
2 **using facial phenotypic descriptors**

3 Hsieh, Tzung-Chien¹; Bar-Haim, Aviram²; Nadav, Guy²; Pantel, Jean Tori^{1,3};
4 Fleischer, Nicole²; Krawitz, Peter^{1,*}

5 1. Institute of Genomic Statistics and Bioinformatics, University of Bonn, Bonn,
6 Germany

7 2. FDNA Inc., Boston Massachusetts, United States

8 3. Charité – Universitätsmedizin Berlin, corporate member of Freie Universität
9 Berlin, Humboldt-Universität zu Berlin, and Berlin Institute of Health, Institute of
10 Medical Genetics and Human Genetics, Berlin, Germany

11 * corresponding author, pkrawitz@uni-bonn.de

12

13 **Abstract**

14 **Introduction:**

15 Recent advances in next-generation phenotyping (NGP) for syndromology, such as
16 DeepGestalt, have learned phenotype representations of multiple disorders by training
17 on thousands of patient photos. However, many Mendelian syndromes are still not
18 represented by existing NGP tools, as only a handful of patients were diagnosed.
19 Moreover, the current architecture for syndrome classification, e.g., in DeepGestalt, is
20 trained “end-to-end,” that is photos of molecularly confirmed cases are presented to
21 the network and a node in the output layer, that will correspond to this syndrome, is
22 maximized in its activity during training. This approach will not be applicable to any
23 syndrome that was not part of the training set, and it cannot explain similarities among
24 patients. Therefore, we propose “GestaltMatch” as an extension of DeepGestalt that
25 utilizes the similarities among patients to identify syndromic patients by their facial
26 gestalt to extend the coverage of NGP tools.

27 **Methods:**

28 We compiled a dataset consisting of 21,400 patients with 1,451 different rare disorders.
29 For each individual, a frontal photo and the molecularly confirmed diagnosis were
30 available. We considered the deep convolutional neural network (DCNN) in
31 DeepGestalt as a composition of a feature encoder and a classifier. The last fully-
32 connected layer in the feature encoder was taken as Facial Phenotypic Descriptor
33 (FPD). We trained the DCNN on the patients’ frontal photos to optimize the FPD and

34 to define a Clinical Face Phenotype Space (CFPS). The similarities among each
35 patient were quantified by cosine distance in CFPS.

36 **Results:**

37 Patients with similar syndromic phenotypes were located in close proximity in the
38 CFPS. Ranking syndromes by distance in CFPS, we first showed that GestaltMatch
39 provides a better generalization of syndromic features than a face recognition model
40 that was only trained on healthy individuals. Moreover, we achieved 87% top-10
41 accuracy in identifying rare Mendelian diseases that were excluded from the training
42 set. We further proved that the distinguishability of syndromic disorders does not
43 correlate with its prevalence.

44 **Conclusions:**

45 GestaltMatch enables matching novel phenotypes and thus complements related
46 molecular approaches.

47

48 **Introduction**

49 Worldwide, rare genetic disorders affect more than 8% of the population. The rarity
50 and diversity of genetic disorders make it time-consuming and challenging for a
51 clinician to achieve a correct diagnosis, which is the so-called “diagnostic odyssey.”¹
52 Craniofacial abnormalities present in 30-40% genetic disorders.² The patients with
53 these syndromic disorders usually show recognizable faces such as Down syndrome
54 and Fragile X syndrome. Hence, the facial manifestation provides a crucial visual hint
55 for a clinician to identify related disorders, which speeds up the diagnostic workup with
56 gene panel or exome sequencing because it helps reduce the search space of
57 candidate genes. However, the ability to recognize these syndromic disorders highly
58 relies on the human expert’s experience. It will be very challenging to make a diagnosis
59 if the clinician has not seen the ultra-rare disorder or novel disease on the patient.
60 Therefore, there is an urgent need for the next-generation phenotyping (NGP) tool to
61 analyze the facial phenotypic information by the aid of a computer.

62 With the rapid development of machine learning and computer vision, a considerable
63 number of NGP tools has emerged for analyzing facial dysmorphology with patients’
64 2D portrait images.³⁻¹⁰ Clinical Face Phenotype Space (CFPS), formed by the facial
65 features extracted from facial images, was proposed to perform syndrome
66 classification on the scale of training on more than 1000 patient photos with eight
67 different syndromes.³ Moreover, face recognition technologies were improved
68 significantly in recent years and were at the core of the deep learning revolution in

69 computer vision. DeepFace¹¹ demonstrated, for the first time, human-level
70 performance in identity verification on the Labeled Faces in the Wild dataset.¹² As a
71 result, the face recognition system trained on CCTV images was utilized to match the
72 patients with one of ten syndromic disorders with intellectual disability.⁶ In addition, the
73 facial recognition model from healthy individuals can also be integrated with the CFPS
74 as a hybrid model, and it was proved to discriminate the facial gestalt on three novel
75 disease-genes.¹⁰ The current state-of-the-art syndrome classification framework
76 DeepGestalt showed record-breaking results for syndrome classification using facial
77 phenotypic cues, achieving 91% top-10 accuracy in identifying the correct syndrome
78 in a test set of 502 images spanning more than 200 syndromes.⁹ DeepGestalt also
79 demonstrated strong separation ability for specific syndromes and subtypes,
80 surpassing human experts' performance. These results demonstrated the power of a
81 community-driven platform to gather patients and collect phenotypes.

82 Although NGP tools have shown the discriminative ability for syndrome classification,
83 they still suffer from the limited data for rare genetic disorders and limited scalability of
84 the model. In Figure 1, the two most well-known studies^{3,9} for multi-syndromes
85 classification focused mainly on the disease-genes with around 50 up to 500
86 pathogenic submissions in ClinVar such as *UBE3A*, *SMC1A*, and *HDAC8* which can
87 be considered as common amongst the rare. The discriminative facial gestalt was
88 identified in *PACS1*, *PPM1D*, and *PHIP*, which moved the border towards the genes
89 with around ten submissions.¹⁰ In addition, two unrelated patients with the same
90 disease-causing mutation in *LEMD2* successfully matched by DeepGestalt syndrome
91 similarity scores.¹³ However, it is still challenging to push the limit to the ultra-rare
92 disease-genes fall in the right tail of the distribution because the NGP approaches
93 require a certain amount of images to learn the facial representation of syndromic
94 disorder.

95 Moreover, the end-to-end offline trained architecture is suboptimal for scaling the
96 model to support new syndromes, to keep the model updated, or to change its original
97 goals. In order to support a new syndrome in DeepGestalt's model, the developer has
98 to go through the six steps described in Supplementary Figure 1. In addition, the model
99 for multi-syndrome classification cannot be used to quantify the similarities among
100 patients that is crucial for clinicians to interpret the patient's phenotype. Therefore, the
101 main limitations of the current approaches are: network architectures that do not scale
102 and that do not allow comparison of single patients.

103 In the nosology of genetic diseases, there has been a discussion about splitters and
104 lumpers for decades.¹⁴ Deep learning approaches cannot only contribute to this

105 dialogue by quantifying distinguishability.¹⁵ The architecture of a well-performing
106 artificial neural network that serves as a classifier for syndromic disorders might reveal
107 something about the complexity of the problem itself. In this work, we consider
108 DeepGestalt as a composition of an image encoder which converts images to a vector
109 of numbers, and a classification head which classifies the encoded vector to soft
110 syndrome probabilities. While the last layer in DeepGestalt consisted of all the
111 syndromes that the network learned to distinguish, we can refer to the layer preceding
112 this last one, as the feature layer.

113 We hypothesize that the new framework, called GestaltMatch, is suitable to
114 overcome these limitations:

- 115 1. Support new syndromes on the fly (without extra training);
- 116 2. Support multiple tasks (e.g., matching patients/syndromes etc.);
- 117 3. Support new explainability approaches (e.g., showing clusters separability in
118 the dataset for different categories);
- 119 4. Be easily customized and allow low maintenance.

120 We show that the features vector created by DeepGestalt encoder can be used as a
121 Facial Phenotypic Descriptor (FPD), which can be further used for syndrome
122 classification and patient clustering. The concept of GestaltMatch is shown in Figure
123 2. Moreover, we show that features created using DeepGestalt encoder are better for
124 matching cases with similar syndromic features, than features extracted from modern
125 models used for face verification and no syndromic phenotype context. Interestingly,
126 we show that our new FPD based framework, named GestaltMatch, has improved
127 scalability for long-tailed syndromes distribution in Figure 1, without the need for
128 retraining. Furthermore, it provides built-in support for patient matching. We show that
129 given a facial image, one can use our system to search patients and syndromes with
130 similar visual phenotypes. Moreover, the similarity between multiple FPDs spans a
131 metric space between syndromes and can be used for finding new phenotypic series
132 or discriminate between affected and non-affected subjects. Our new system is a
133 natural extension to DeepGestalt and can help to develop new visual phenotype
134 matching applications.

135 **Method**

136 **Datasets**

137 We collected the images of subjects with clinically or molecularly confirmed diagnoses
138 from Face2Gene database. The images with poor quality or duplicated images were

139 removed from the dataset. After removing the problematic images, the dataset
140 consisted of 33,434 images and 21,400 subjects of 1451 syndromes in total.

141 GestaltMatch aims to evaluate syndromes with different properties. We separated the
142 syndromes by the number of subjects in each syndrome and whether they were
143 learned by the DeepGestalt model. The overview of how the dataset was divided is
144 shown in Supplementary Figure 2. The current DeepGestalt approach can only learn
145 the syndromes which have more than six subjects. Hence, based on this threshold, we
146 first separated the syndromes into frequent and rare syndromes. We denoted rare
147 syndromes as target syndromes because these are the syndromes on which this study
148 targets. However, not all frequent syndromes can be modeled by DeepGestalt. Some
149 of them might have no dysmorphic features, so DeepGestalt cannot learn their facial
150 representation. We denoted these syndromes as non-distinct, whereas the syndromes
151 supported by DeepGestalt as distinct. The distinct syndromes were used for validating
152 syndrome prediction and the separation ability of subtypes of a phenotypic series
153 because these syndromes were known to have facial dysmorphic features, and the
154 facial features were well recognized by DeepGestalt encoder. For target syndromes,
155 we aim to prove that GestaltMatch is able to predict the syndrome even if only a few
156 subjects are publicly available. It is noteworthy that currently, for more than half of all
157 known disease-genes, less than ten cases with pathogenic mutations have been
158 submitted to ClinVar (Figure 1). By the type of syndromes, we split the entire dataset
159 into three datasets: distinct, non-distinct, and target syndromes, and they contained
160 301, 265, and 885 syndromes, respectively. Non-distinct and target syndromes are not
161 yet applicable to DeepGestalt.

162 We further sampled each dataset into a gallery and test set. The gallery is a set of
163 subjects we intend to match, given a subject from the test set. First of all, 1422 subjects
164 in distinct and non-distinct datasets were kept out of the training set as a blind set for
165 validating the DeepGestalt training. The subjects in the blind set were assigned to
166 either distinct test set or non-distinct test set based on the type of syndromes, and the
167 subjects not in the blind set were assigned to the gallery of the corresponding dataset.
168 For the target dataset, we performed 10-fold cross-validation. 90% and 10% of
169 subjects were assigned to the gallery and test set, respectively.

170 However, if we only performed the matching within the same dataset, it will not be the
171 real-world scenario. The galleries of three datasets were later combined as a unified
172 gallery, and we try to find the matched patients in the unified gallery. We called the
173 gallery of each dataset as a partial gallery. It is used for the performance comparison
174 between the DeepGestalt model and GestaltMatch on distinct syndromes because

175 DeepGestalt only predicts distinct syndromes, so we should only use the partial gallery
176 of the distinct set as the gallery.

177 **DeepGestalt encoder**

178 The preprocessing pipeline of DeepGestalt includes points detection, facial alignment
179 (frontalization), and facial regions cropping. During inference, every facial region crop
180 is forward passed through a deep convolutional network (CNN), and finally, the results
181 for all of the image regions are aggregated to the final prediction for the input face
182 image. DeepGestalt network consists of ten convolutional layers with batch
183 normalization (BN) and ReLU for embedding the input features. After every Conv-BN-
184 ReLU layer, a max pooling layer is applied for reducing the spatial size while increasing
185 the semantic representation. The classifier part of the network consists of a fully
186 connected linear layer with dropout (0.5). In this work, we considered DeepGestalt
187 architecture as an encoder-classification composition, pipelined during inference. We
188 chose the last fully connected layer before the softmax classification as the facial
189 feature representation, resulting in a vector of size 320. Our first hypothesis is that
190 images with the same molecularly diagnosed syndromes or phenotypic series, which
191 also share similar phenotypes, can be encoded to similar feature vectors, under some
192 set of metrics.

193 Moreover, we claim that the specific design choice of DeepGestalt of using a
194 predefined, offline trained, linear classifier, can be replaced by other classification
195 “heads,” for example, k -Nearest Neighbors using cosine distance or a Random Forest.
196 Interestingly, we found that the data used during the FPD encoder training is essential
197 to generalize unseen syndromes, subjects, and the space represented by the FPD
198 encoder.

199 **Descriptor projection - Clinical Face Phenotype Space**

200 Each image was encoded by the DeepGestalt encoder and resulted in a 320-
201 dimensional facial phenotypic descriptor. These facial phenotypic descriptors were
202 further used to form a 320-dimensional space which is called Clinical Face Phenotype
203 Space (CFPS), and each image is a point located in CFPS, as shown in Figure 2. The
204 similarity between the two images is quantified by the cosine distance between them
205 in CFPS. The smaller the distance is, the higher similarity between two images is.
206 Therefore, the clusters of subjects in CFPS can represent the similarities among the
207 different disorders or show the substructure under a phenotypic series.

208 **Evaluation**

209 To evaluate GestaltMatch, we take the images in the test set as input and position
210 them in the CFPS that is defined by the images of the gallery. We calculated the cosine
211 distance between each of the test set images to all the gallery images, and benchmark
212 the performance by top- k accuracy. For each test image, if an image from another
213 subject with the same disorder in the gallery is among the top- k nearest neighbors, we
214 call it a top- k match. We further compare the accuracy of each syndrome in distinct,
215 non-distinct, and target syndrome subsets to investigate whether GestaltMatch can
216 extend DeepGestalt to support more syndromes.

217 **Results**

218 **Comparing DeepGestalt and face recognition encoders**

219 We first investigated the importance of using a syndromic features encoder rather than
220 a normal facial features encoder. We compared FPDs created by DeepGestalt
221 encoder to another encoder with the same architecture, trained on the CASIA-
222 WebFace¹⁶ recognition task. We then trained these two encoders and encoded all
223 images by these two encoders separately.

- 224 • Enc-DeepGestalt encoder, trained on the gallery of 301 distinct syndromes.
- 225 • Enc-CASIA encoder, trained on the CASIA-WebFace dataset, with the same
226 architecture of DeepGestalt.

227 We evaluated the performance by testing distinct and target test sets on the unified
228 gallery. Table 1 shows the superiority of the features created by DeepGestalt in the
229 matching performance, which emphasizes the importance of training the encoder on
230 data with phenotypic cues. The features created by DeepGestalt improves the top-10
231 accuracy by 30% for the distinct category. Further, the top-10 accuracy was improved
232 by 43% for the target syndromes, which contains a different, mutually exclusive list of
233 syndromes. These results suggest that the features encoded by DeepGestalt are a
234 better fit for the task of syndromes classification than the features encoded by the
235 modern face recognition model. Moreover, DeepGestalt’s FPD provides a better
236 generalization than the FPD encoded by the modern face recognition model for unseen
237 target syndromes.

238 **Comparing distinct and non-distinct FPDs**

239 In order to demonstrate the separability of syndromes with facial dysmorphism, we
240 applied t -SNE¹⁷ to project 4353 images of ten distinct syndromes with the largest
241 number of subjects and 872 images of ten non-distinct syndromes to two-dimensional
242 space, and we further calculated Silhouette index¹⁸ for both of two datasets. Autism

243 syndrome has 1171 images which is the largest non-distinct syndrome. We did not
244 take Autism into this analysis because it leads to an extreme imbalance of the number
245 of subjects to the other non-distinct syndromes. As shown in Supplementary Figure 3,
246 the FPDs of distinct syndrome show ten clear clusters of subjects. However, when
247 applying *t*-SNE projection on subjects of non-distinct syndromes, no clear clusters are
248 created. Besides, the Silhouette index of distinct syndromes is 0.07, which is higher
249 than the index of non-distinct syndromes, which is -0.01. The negative Silhouette index
250 of non-distinct syndromes indicates the poor separation of different syndromes. The
251 results show the evidence for the phenotypic information encoded in the FPDs created
252 by DeepGestalt.

253 **GestaltMatch on unseen dysmorphic syndromes**

254 For the purpose of proving GestaltMatch can match the patients with a novel syndrome
255 unseen to the encoder and to better understand the important characteristics of the
256 training dataset, we trained four different encoders for comparison. We sampled 21228
257 images of 13872 subjects with 279 known dysmorphic syndromes. The four encoder
258 variants Enc-1 to Enc-4 are:

- 259 • Enc-1, trained on 90% of the 279 syndromes' subjects;
- 260 • Enc-2, trained on 90% of the 239 smallest syndromes' subjects;
- 261 • Enc-3, trained on 90% of the 239 largest syndromes' subjects;
- 262 • Enc-4, trained on 90% of 239 random syndromes' subjects.

263 For each model, we used the remaining 10% of the subjects, sampled across the
264 syndromes in the training set, as a validation set. Moreover, we used the remaining 40
265 syndromes, of encoders 2-4 (the eliminated syndromes for each encoder) as an
266 external test set of unseen distinct syndromes, denoted by Test-Large, Test-Small,
267 and Test-Random, respectively. For example, the 40 syndromes in Test-Large are the
268 largest 40 syndromes in 279 distinct syndromes, complementing the 239 syndromes
269 trained in Enc-2.

270 To evaluate FPDs generalization ability of each encoder on unseen syndromes, we
271 compared each of three encoders (Enc-2, Enc-3, and Enc-4), trained on a subset of
272 239 syndromes, to Enc-1. We used GestaltMatch to estimate the similarity between
273 the test images to the gallery images, with cosine distance. The results are shown in

274 Table 2. Enc-1 outperformed Enc-2 when testing on Test-Large, the top-10 accuracy
275 dropped from 85.28% to 78.49%. The poor performance of Enc-2 could be due to
276 losing too much training data because Test-Large contained the largest 40 distinct
277 syndromes consisting of 12429 images, which is more than half of the total images.

278 However, Enc-3 and Enc-4 showed comparable results to Enc-1 on Test-Small and
279 Test-Random, respectively. The top-10 accuracy of Enc-4 was even slightly higher
280 than Enc-1 when testing on Test-Random. It means that the encoder without the 40
281 smallest or random distinct syndromes, leads to comparable performance on these 40
282 syndromes with an encoder trained with these 40 syndromes. Therefore, the results
283 proved that GestaltMatch could generalize the facial dysmorphic features well on
284 unseen syndromes, which means we are able to support new syndrome without
285 retraining the model.

286 **Target syndromes matching accuracy**

287 We defined a syndrome as a target syndrome if it has less than seven subjects in our
288 dataset. To understand the potential of matching target syndromes, we trained an
289 encoder on 2215 images of 526 target syndromes, which have more than three
290 subjects, and less than seven subjects, denoted by Enc-Target. We then compared
291 Enc-Target to Enc-DeepGestalt trained on the 301 distinct syndromes from the
292 previous section. Results in

293

294 Table 3 show that Enc-Target with a softmax classifier provides the best results, which
295 means it learned important phenotypic features. However, in GestaltMatch only, Enc-
296 DeepGestalt, which trained on distinct syndromes and did not see any of the target
297 syndromes during its training, showed very similar results compared to Enc-Target.
298 Although the results showed that cosine distance is inferior to a trained softmax
299 classifier, encoders trained on distinct syndromes provide a similar accuracy on the
300 unseen subject of target syndromes, compared to encoders trained on these target
301 syndromes. Therefore, GestaltMatch is a more suitable choice for target syndromes
302 because it achieved comparable performance to the encoder train on target syndrome,
303 that means we could save resources for retraining the encoder. Moreover, training the
304 model on both distinct and target syndromes, which have very few high-quality photos,
305 might lead to poor performance due to the extremely imbalanced training dataset.

306 **Correlation between prevalence and accuracy**

307 Training an end-to-end network for classifying faces to syndromes such as
308 DeepGestalt requires many subjects for each of the supported syndromes. Since this
309 minimum subject requirement is no longer a must for GestaltMatch, we were interested
310 in whether the matching accuracy of a syndrome correlates with its prevalence. We
311 used Enc-2 from the previous section, trained on the 239 syndromes out of the full list
312 of 279 syndromes. To remove the confounding effect from prevalence, we randomly

313 down-sampled each of the 40 avoided syndromes by selecting five subjects to the
314 gallery and one subject to the test set. This experiment is repeated 1000 times. Figure
315 3 shows the average top-10 accuracy and the prevalence by Orphanet of each
316 syndrome. We can see that the top-10 accuracy does not correlate with the prevalence.
317 Several syndromes with low prevalence still perform very well. We can further consider
318 the accuracy as the distinguishability of the syndrome. Therefore, as distinguishability
319 does not depend on the prevalence, GestaltMatch can extend DeepGestalt to cover
320 the ultra-rare disorders with a high distinguishability.

321 **Hierarchical clustering**

322 Phenotypic series is defined as a heterogeneous set of genetic disorders sharing
323 similar phenotypes. We were interested in testing the visual clusters created with a
324 two-dimensional projection of FPDs. We sampled subjects from subtypes of four large
325 phenotypic series in our database: Noonan syndrome, Cornelia De Lange Syndrome
326 (CDLS), Kabuki syndrome, and Mucopolysaccharidosis (MPS). As demonstrated in
327 Supplementary Figure 4, using *t*-SNE projection on the FPDs of 743 subjects, sampled
328 from the four phenotypic series, resulted in highly separable four clusters composed
329 of the different subtypes of each phenotypic series. This result is a piece of evidence
330 for the phenotypic features encoded in the FPDs.

331 **Dysmorphism estimation**

332 We were interested in testing GestaltMatch separation ability between FPDs of
333 affected subjects with a dysmorphic genetic disorder and non-affected subjects
334 (without a known genetic disorder). We used *t*-SNE projections in two different formats:

- 335 1. We sampled 1000 faces of healthy individuals and added to the ten largest
336 app-valid syndromes (4346 subjects) projection;
- 337 2. We sampled 1000 faces of healthy individuals and 1000 faces of non-healthy
338 individuals, evenly across the ten largest app-valid syndromes.

339 In the first experiment, we projected into eleven classes, while in the second
340 experiment, we used a binary classification into two clusters. Results in Supplementary
341 Figure 5 and Supplementary Figure 6 show that in both cases, non-affected subjects
342 create reasonably separable clusters, emphasizing the syndromic context encoded
343 in the FPDs by GestaltMatch.

344 **Discussion**

345 **Syndrome matching**

346 As described in the evaluation section, the GestaltMatch framework can be used to
347 match syndromes with an input image. The difference from DeepGestalt, lies in the
348 ability to match unseen syndromes (no patient with these syndromes included in the
349 encoder's training set). GestaltMatch framework also allows us to abstract away the
350 encoding of a dataset from the classification task, and thus support multiple targets
351 within the same evaluation. For example, one can evaluate both phenotypic series and
352 subtypes levels within a single inference, or get the most similar patients for each of
353 the matched syndromes with a minor computational cost which is only a few seconds.

354 GestaltMatch framework computes the similarity between each of the test set images
355 to the entire dataset of images. The similarity can be computed using different metrics,
356 for example, cosine or euclidean distance. Then the results are aggregated according
357 to the chosen configuration. For example, image similarity can be aggregated at the
358 patient level or in the syndrome level. Furthermore, filtering according to different
359 dataset parameters (such as ethnicity, number of affected genes, and age), can be
360 done to customize the evaluation further.

361 **Patient matching**

362 Matching patients with high similarities of facial dysmorphic features is one of the most
363 important applications of GestaltMatch. Finding the second patient is always a
364 challenging problem for most of the physicians when analyzing the novel or extremely
365 rare Mendelian disorders. There are several online platforms such as Gene Matcher,¹⁹
366 MyGene2, and Exchange Maker,²⁰ which allow physicians to look for similar patients
367 by uploading phenotypic data, such as HPO terms or genomic information. These
368 platforms have already matched thousands of patients in the past few years. However,
369 the automated facial matching technology was not included in any of these platforms
370 yet, although the facial phenotypes are crucial information for physicians to determine
371 whether two patients have similar disorders or not. Therefore, there is an urgent need
372 to support the patient matching approach by analyzing the facial images to facilitate
373 the matching procedure.

374 As a proof of concept, we have matched two unrelated patients from different countries,
375 with the same novel disease successfully by gestalt matching approach.¹³ They both
376 shared similar progeria-like features, and later the same *de novo* disease-causing
377 mutation in the *LEMD2* gene was identified by the diagnostic workup. Although further
378 analysis with more unrelated patients is needed to be done, the GestaltMatch
379 approach could be a promising patient match application. Moreover, this approach can
380 be integrated with the other matching platforms to enhance the matching ability to

381 reduce the amount of time further when looking for the second patient with the same
382 rare disorder.

383 **Ethnic bias**

384 Ethnic bias could influence the performance of GestaltMatch dramatically, especially
385 in target syndromes, because some patients with the same target syndromes were
386 from the same paper. Moreover, most of those patients with extremely rare diseases
387 in the same publication are usually from the same family. It is hard to tell whether they
388 were matched by the dysmorphic features or the ethnic similarity. Therefore, testing
389 the matching of subjects of distinct syndromes with different ethnic backgrounds by a
390 statistical setting to assess the influence of ethnic bias is needed in the future
391 experiment.

392 **The future of phenotype representation**

393 Using semantic descriptors for similarity estimation is common in many areas of
394 artificial intelligence, such as face recognition, text understanding, visual tracking,
395 speech recognition, and more. In recent years, converting structured data into
396 semantic vectors is becoming common for non-visual phenotype matching as well, for
397 example, in HPO2VEC²¹ and NODE2VEC.²²

398 Moreover, to improve the matching accuracy, input signals can be sourced from
399 different modalities. For example, DeepGestalt uses different regions of a patient face
400 and aggregates the classification result of each region. Moreover, in PEDIA,²³
401 semantic and visual phenotypic cues are aggregated to improve the prioritization of
402 variant analysis.

403 Since semantic descriptors share the same format, one can aggregate these
404 descriptors from different sources, to allow multimodal signals contribution to the final
405 accuracy. Due to the generic structure of the GestaltMatch framework and the
406 abstractions used for encoding datasets, future work can extend the GestaltMatch
407 framework to support different input types such as text, speech, video, or other sources
408 of medical imaging, to improve classification accuracy.

409 **Designing a unified classification approach**

410 One of the main challenges of productionizing the GestaltMatch technology lies in the
411 ability to aggregate different categories. As shown in unseen syndromes analysis, an
412 internal bias in the encoder's dataset can deteriorate the matching performance for
413 both seen and unseen syndromes. Moreover, training a softmax classifier (as in
414 DeepGestalt) provides better accuracy than a naive cosine distance over FPDs. The

415 question raised from these insights is - how to use GestaltMatch for supporting all types
416 of syndromes? Accurately, future work will test whether it is better to combine
417 GestaltMatch classification (for unseen or target syndromes) and DeepGestalt (for
418 distinct syndromes) in a hybrid manner or use a single model to directly classify an
419 image to all syndromes (using GestaltMatch or DeepGestalt).

420 **Conclusion**

421 GestaltMatch can match syndromes with facial dysmorphism in the CFPS and can be
422 treated as an extension of DeepGestalt to cover the syndromes which are not
423 supported in DeepGestalt model. Moreover, the sub-structure under a phenotypic
424 series or novel diseases can be explored by the clustering of subjects in CFPS.
425 Eventually, matching patients is one of the most important applications of GestaltMatch.
426 It could be integrated into other online matching platforms such as MatchMaker
427 Exchange or MyGene2 further to accelerate the matching process of unknown
428 diagnosed patients and explore novel phenotype-genotype correlation.

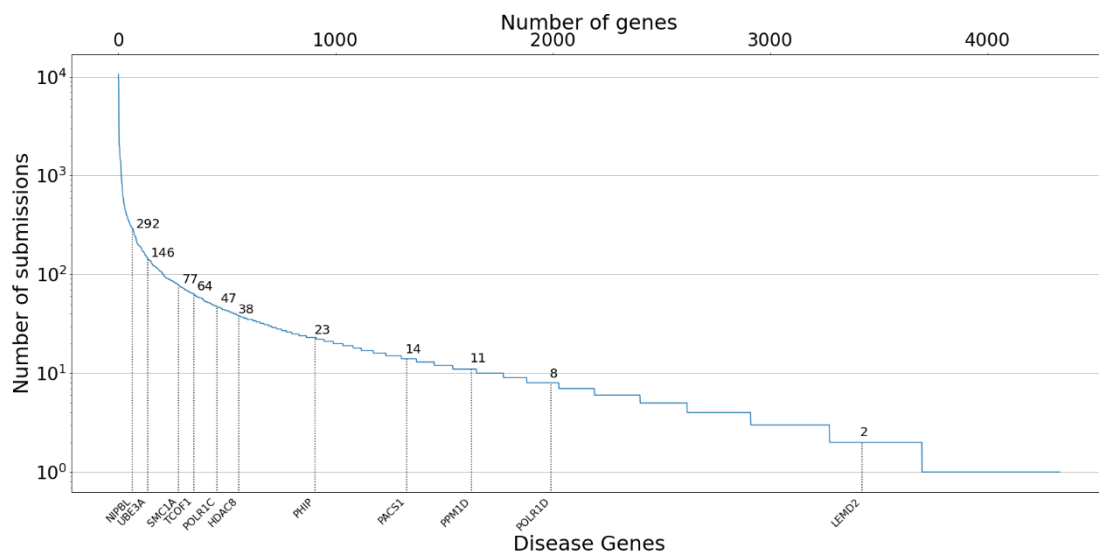
429 **Reference**

- 430 1. Baird, P. A., Anderson, T. W., Newcombe, H. B. &Lowry, R. B. Genetic
431 disorders in children and young adults: A population study. *Am. J. Hum.*
432 *Genet.* **42**, 677–693 (1988).
- 433 2. Hart, T. &Hart, P. Genetic studies of craniofacial anomalies: clinical
434 implications and applications. *Orthod. Craniofac. Res.* **12**, 212–220 (2009).
- 435 3. Ferry, Q. *et al.* Diagnostically relevant facial gestalt information from ordinary
436 photos. 1–22 (2014). doi:10.7554/eLife.02020
- 437 4. Cerrolaza, J. J. *et al.* Identification of dysmorphic syndromes using landmark-
438 specific local texture descriptors. *2016 IEEE 13th International Symposium on*
439 *Biomedical Imaging (ISBI)* 1080–1083 (2016). doi:10.1109/ISBI.2016.7493453
- 440 5. Wang, K. &Luo, J. Detecting Visually Observable Disease Symptoms from
441 Faces. *EURASIP J. Bioinform. Syst. Biol.* **2016**, 13 (2016).
- 442 6. Dudding-Byth, T. *et al.* Computer face-matching technology using two-
443 dimensional photographs accurately matches the facial gestalt of unrelated
444 individuals with the same syndromic form of intellectual disability. *BMC*
445 *Biotechnol.* **17**, 1–9 (2017).
- 446 7. Shukla, P., Gupta, T., Saini, A., Singh, P. &Balasubramanian, R. A Deep
447 Learning Frame-Work for Recognizing Developmental Disorders. *2017 IEEE*
448 *Winter Conference on Applications of Computer Vision (WACV)* 705–714

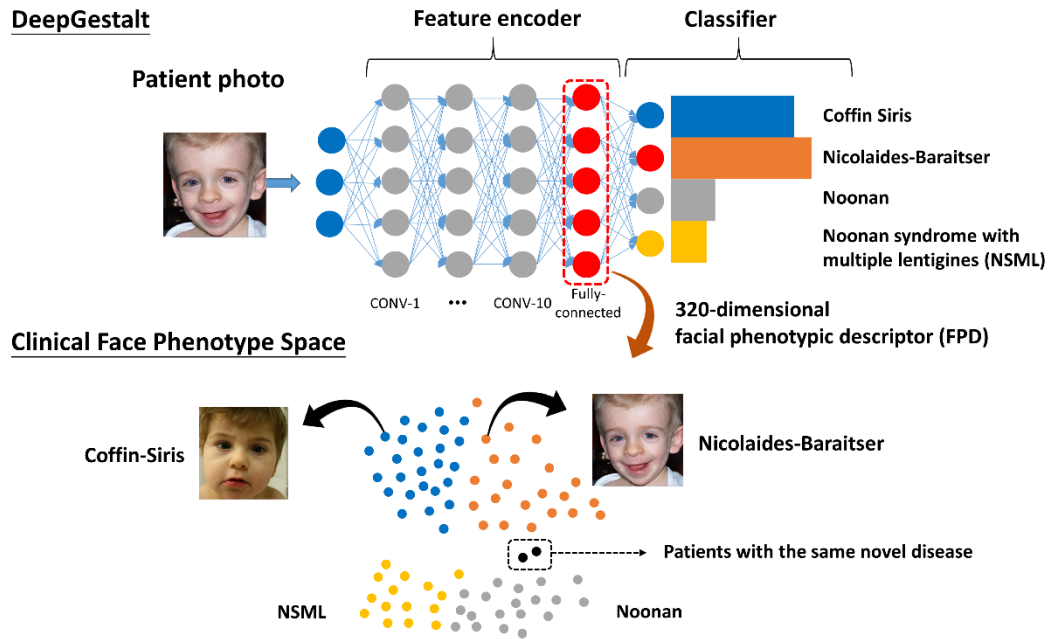
- 449 (2017). doi:10.1109/WACV.2017.84
- 450 8. Liehr, T. *et al.* Next generation phenotyping in Emanuel and Pallister-Killian
451 syndrome using computer-aided facial dysmorphology analysis of 2D photos.
452 *Clin. Genet.* **93**, 378–381 (2018).
- 453 9. Gurovich, Y. *et al.* Identifying facial phenotypes of genetic disorders using
454 deep learning. *Nature Medicine* **25**, 60–64 (2019).
- 455 10. van derDonk, R. *et al.* Next-generation phenotyping using computer vision
456 algorithms in rare genomic neurodevelopmental disorders. *Genet. Med.* **21**,
457 1719–1725 (2019).
- 458 11. Taigman, Y., Yang, M., Ranzato, M. & Wolf, L. DeepFace: Closing the gap to
459 human-level performance in face verification. in *Proceedings of the IEEE*
460 *Computer Society Conference on Computer Vision and Pattern Recognition*
461 1701–1708 (IEEE Computer Society, 2014). doi:10.1109/CVPR.2014.220
- 462 12. Learned-Miller, E., Huang, G. B., RoyChowdhury, A., Li, H. & Hua, G. Labeled
463 faces in the wild: A survey. *Adv. Face Detect. Facial Image Anal.* 189–248
464 (2016). doi:10.1007/978-3-319-25958-1_8
- 465 13. Marbach, F. *et al.* The Discovery of a LEMD2-Associated Nuclear Envelopathy
466 with Early Progeroid Appearance Suggests Advanced Applications for AI-
467 Driven Facial Phenotyping. *Am. J. Hum. Genet.* **104**, 749–757 (2019).
- 468 14. McKusick, V. A. On lumpers and splitters, or the nosology of genetic disease.
469 *Perspect. Biol. Med.* **12**, 298–312 (1969).
- 470 15. Knaus, A. *et al.* Characterization of glycosylphosphatidylinositol biosynthesis
471 defects by clinical features, flow cytometry, and automated image analysis.
472 *Genome Med.* **10**, 3 (2018).
- 473 16. Yi, D., Lei, Z., Liao, S. & Li, S. Z. Learning Face Representation from Scratch.
474 (2014).
- 475 17. Van DerMaaten, L. & Hinton, G. *Visualizing Data using t-SNE.* **9**, (2008).
- 476 18. Rousseeuw, P. J. Silhouettes: a graphical aid to the interpretation and
477 validation of cluster analysis. *J. Comput. Appl. Math.* **20**, 53–65 (1987).
- 478 19. Sobreira, N., Schiettecatte, F., Valle, D. & Hamosh, A. GeneMatcher: A
479 Matching Tool for Connecting Investigators with an Interest in the Same Gene.
480 *Hum. Mutat.* **36**, 928–930 (2015).
- 481 20. Philippakis, A. A. *et al.* The Matchmaker Exchange: A Platform for Rare
482 Disease Gene Discovery. *Hum. Mutat.* **36**, 915–921 (2015).

- 483 21. Shen, F. *et al.* HPO2Vec+: Leveraging heterogeneous knowledge resources to
 484 enrich node embeddings for the Human Phenotype Ontology. *J. Biomed.*
 485 *Inform.* **96**, 103246 (2019).
- 486 22. Grover, A. & Leskovec, J. Node2vec: Scalable feature learning for networks. in
 487 *Proceedings of the ACM SIGKDD International Conference on Knowledge*
 488 *Discovery and Data Mining 13-17-Aug*, 855–864 (Association for Computing
 489 Machinery, 2016).
- 490 23. Hsieh, T. C. *et al.* PEDIA: prioritization of exome data by image analysis.
 491 *Genet. Med.* **21**, 2807–2814 (2019).

492 Figures and tables

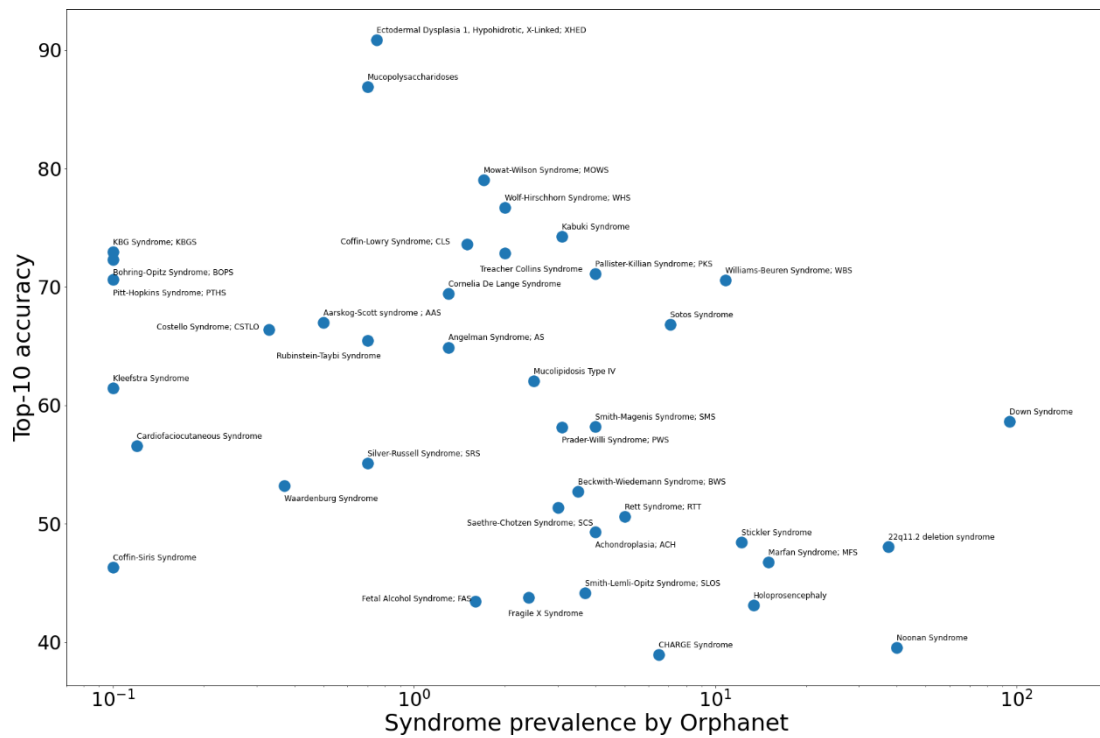


493
 494 **Figure 1: The distribution of the number of pathogenic submissions of each gene in**
 495 **ClinVar (April, 2020).** The lower x-axis shows the disease genes, and the upper x-axis is the
 496 number of genes cumulative from zero on the left. Y-axis is the number of pathogenic
 497 submissions in ClinVar for the respective gene. The most two well-known multi-syndromes
 498 classification studies^{3,9} mainly focused on the syndromes with relative common in the rare
 499 disorders such as Angelman syndrome (*UBE3A*), Cornelia de Lange syndrome (*NIPBL*,
 500 *SMC1A*, *HDAC8*), and Treacher Collins syndrome (*TCOF1*, *POLR1D*). The three novel
 501 disease-genes (*PACS1*, *PPM1D*, and *PHIP*,) which is proved to show discriminative facial
 502 gestalt,¹⁰ were relatively rare compared to the previous studies. Later, the new disease related
 503 to *LEMD2* was found by two matching patients with similar facial phenotype.¹³ *LEMD2* even
 504 only had two submissions so far. It shows that NGP approaches keep pushing the limit to more
 505 ultra-rare diseases on the right tail. However, 59% (2562 out of 4335) of disease genes with
 506 less than ten pathogenic submissions in ClinVar. The limited patients of rare disorders are a
 507 challenge to the current NGP approach since it requires a certain number of images to learn
 508 the facial representation of a disorder.



509

510 **Figure 2: Concept of GestaltMatch.** The DeepGestalt was trained on 301 distinct syndromes
 511 to learn the facial dysmorphic features. The last fully-connected layer in the feature encoder is
 512 taken as Facial Phenotypic Descriptor (FPD) and can be used to form a Clinical Face
 513 Phenotype Space (CFPS). In this space, the distance between each patient can be considered
 514 as the similarities of facial phenotypic features, which can be further used for syndrome
 515 classification or clustering patients with unknown diagnosis.



516

517 **Figure 3: Correlation between syndrome prevalence and average top-10 accuracy.** X-axis

518 is the birth prevalence by Orphanet, and the unit is 1 in 100,000. For each syndrome, we
 519 randomly selected five subjects to the gallery and one subject to the test set to remove the
 520 confounding effect from prevalence. We further performed the classification on these 40
 521 syndromes for 1000 times. Y-axis is the average top-10 accuracy of the experiments with 1000
 522 times. From this figure, we can see that the top-10 accuracy does not correlate with disease
 523 prevalence.

524 **Table 1: Performance comparison of DeepGestalt and CASIA encoder on distinct, non-**
 525 **distinct and target test set.** Enc-DeepGestalt and Enc-CASIA have the same architecture.
 526 Enc-DeepGestalt was initiated with the CASIA-WebFace and further fine-tuned on patients'
 527 photos. Enc-DeepGestalt outperformed Enc-CASIA on distinct and target syndromes. It shows
 528 the importance of fine-tuning on patients' photos for learning facial dysmorphic features.

Test set	Model	Syndromes		Top 1	Top 5	Top 10	Top 30
		Gallery	Test				
Distinct	Enc-DeepGestalt	1451	168	33.46%	56.11%	66.73%	80.54%
Distinct	Enc-CASIA	1451	168	20.10%	41.49%	51.33%	70.23%
Non-distinct	Enc-DeepGestalt	1451	75	6.44%	11.85%	15.79%	27.99%
Non-distinct	Enc-CASIA	1451	75	7.30%	11.16%	14.25%	20.26%
Target	Enc-DeepGestalt	1451	885	6.84%	12.97%	16.03%	21.69%
Target	Enc-CASIA	1451	885	4.67%	8.49%	11.21%	15.41%

529

530 **Table 2: Results of unseen syndromes classification with four encoders.** Each pair of
 531 results below shows the comparison between training without the syndromes in the test set and
 532 with them. For example, Enc-1 was trained on 279 distinct syndromes, and Enc-2 was trained
 533 on the 239 distinct syndromes. The 40 syndromes in Test-Large are the unseen syndromes to
 534 Enc-2. When testing on Test-Random, Enc-4 shows the comparable results to Enc-1.

Test set	Model	Images		Syndromes	Top 1	Top 5	Top 10	Top 30
		Gallery	Test					
Test-Large	Enc-2	12429	1311	40	32.57%	65.45%	78.49%	97.79%
Test-Large	Enc-1	12429	1311	40	44.85%	74.07%	85.28%	98.32%
Test-Small	Enc-3	532	87	40	37.93%	73.56%	82.76%	96.55%
Test-Small	Enc-1	532	87	40	44.83%	67.82%	85.06%	97.70%
Test-Random	Enc-4	4025	430	40	47.44%	77.44%	87.44%	99.07%
Test-Random	Enc-1	4025	430	40	53.02%	77.67%	86.74%	99.07%

535

536

537

538 **Table 3: Comparison of different models for matching target syndromes.** Enc-
 539 DeepGestalt is the encoder trained on 301 distinct syndromes, and Enc-Target is the encoder
 540 trained on 526 target syndromes. The last row used DeepGestalt method, which is the softmax
 541 in DeepGestalt model to predict the syndrome, so it did not use the gallery.

Model	Method	Images		Syndromes	Top 1	Top 5	Top 10	Top 30
		Gallery	Test					
Enc-DeepGestalt	GestaltMatch	2215	749	526	14.81%	23.98%	29.57%	41.84%
Enc-Target	GestaltMatch	2215	749	526	14.55%	24.70%	30.04%	42.59%
Enc-Target	DeepGestalt	-	749	526	17.35%	26.56%	32.84%	44.30%

542