

# 1 **GestaltMatcher: Overcoming the limits of rare disease** 2 **matching using facial phenotypic descriptors**

3 Tzung-Chien Hsieh<sup>1,+</sup>, Aviram Bar-Haim<sup>2,+</sup>, Shahida Moosa<sup>3</sup>, Nadja Ehmke<sup>4</sup>, Karen  
4 W. Gripp<sup>5</sup>, Jean Tori Pantel<sup>1,4</sup>, Magdalena Danyel<sup>4,6</sup>, Martin Atta Mensah<sup>4,7</sup>, Denise  
5 Horn<sup>4</sup>, Stanislav Rosnev<sup>4</sup>, Nicole Fleischer<sup>2</sup>, Guilherme Bonini<sup>2</sup>, Alexander Hustinx<sup>1</sup>,  
6 Alexander Schmid<sup>1</sup>, Alexej Knaus<sup>1</sup>, Behnam Javanmardi<sup>1</sup>, Hannah Klinkhammer<sup>1,8</sup>,  
7 Hellen Lesmann<sup>1</sup>, Sugirthan Sivalingam<sup>1,8,9</sup>, Tom Kamphans<sup>10</sup>, Wolfgang  
8 Meiswinkel<sup>10</sup>, Frédéric Ebstein<sup>11</sup>, Elke Krüger<sup>11</sup>, Sébastien Küry<sup>12,13</sup>, Stéphane  
9 Bézieau<sup>12,13</sup>, Axel Schmidt<sup>14</sup>, Sophia Peters<sup>14</sup>, Hartmut Engels<sup>14</sup>, Elisabeth Mangold<sup>14</sup>,  
10 Martina Kreiß<sup>14</sup>, Kirsten Cremer<sup>14</sup>, Claudia Perne<sup>14</sup>, Regina C. Betz<sup>14</sup>, Tim  
11 Bender<sup>14,15</sup>, Kathrin Grundmann-Hauser<sup>16</sup>, Tobias B. Haack<sup>16</sup>, Matias Wagner<sup>17,18</sup>,  
12 Theresa Brunet<sup>17</sup>, Heidi Beate Bentzen<sup>19</sup>, Luisa Averdunk<sup>20</sup>, Kimberly Christine  
13 Coetzer<sup>3</sup>, Gholson J. Lyon<sup>21,22</sup>, Malte Spielmann<sup>23</sup>, Christian Schaaf<sup>24</sup>, Stefan  
14 Mundlos<sup>4</sup>, Markus M. Nöthen<sup>14</sup>, Peter Krawitz<sup>1,\*</sup>

15

16 <sup>1</sup>Institute for Genomic Statistics and Bioinformatics, University Hospital Bonn,  
17 Rheinische Friedrich-Wilhelms-Universität Bonn, Bonn, Germany;

18 <sup>2</sup>FDNA Inc., Boston, MA, United States;

19 <sup>3</sup>Division of Molecular Biology and Human Genetics, Stellenbosch University and  
20 Medical Genetics, Tygerberg Hospital, Tygerberg, South Africa;

21 <sup>4</sup>Institute of Medical Genetics and Human Genetics, Charité-Universitätsmedizin Berlin,  
22 Humboldt-Universität zu Berlin and Berlin Institute of Health, Berlin, Germany;

23 <sup>5</sup>A.I. DuPont Hospital for Children/Nemours, Wilmington, DE, USA;

24 <sup>6</sup>Berlin Center for Rare Diseases, Charité-Universitätsmedizin Berlin, Humboldt-  
25 Universität zu Berlin and Berlin Institute of Health, Berlin, Germany;

26 <sup>7</sup>Berlin Institute of Health (BIH), Berlin, Germany;

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

- 27 <sup>8</sup>Institute for Medical Biometry, Informatics and Epidemiology, Medical  
28 Faculty, University of Bonn, Bonn, Germany;
- 29 <sup>9</sup>Core Unit for Bioinformatics Data Analysis, Medical Faculty, University  
30 of Bonn, Bonn, Germany;
- 31 <sup>10</sup>GeneTalk, Bonn, Germany;
- 32 <sup>11</sup>Institut für Medizinische Biochemie und Molekularbiologie (IMBM),  
33 Universitätsmedizin Greifswald, Greifswald, Germany;
- 34 <sup>12</sup>CHU Nantes, Service de Génétique Médicale, Nantes, France;
- 35 <sup>13</sup>Institut du Thorax, INSERM, CNRS, Université de Nantes, Nantes, France;
- 36 <sup>14</sup>Institute of Human Genetics, University of Bonn, Medical Faculty & University  
37 Hospital Bonn, Bonn, Germany;
- 38 <sup>15</sup>Center for Rare Diseases Bonn, University Hospital Bonn, Bonn, Germany;
- 39 <sup>16</sup>Institute of Medical Genetics and Applied Genomics, University of Tübingen,  
40 Tübingen, Germany;
- 41 <sup>17</sup>Institute of Human Genetics, School of Medicine, Technical University Munich,  
42 Munich, Germany;
- 43 <sup>18</sup>Institute of Neurogenomics, Helmholtz Zentrum München GmbH, German Research  
44 Center for Environmental Health, Neuherberg, Germany;
- 45 <sup>19</sup>Norwegian Research Center for Computers and Law, Faculty of Law, University of  
46 Oslo, Oslo, Norway;
- 47 <sup>20</sup>Institute of Human Genetics and Department of Pediatrics, Medical Faculty, Heinrich  
48 Heine University, Düsseldorf, Germany;
- 49 <sup>21</sup>Department of Human Genetics and George A. Jervis Clinic, NYS Institute for Basic  
50 Research in Developmental Disabilities, Staten Island NY 10314, USA;
- 51 <sup>22</sup>Biology PhD Program, The Graduate Center, The City University of New York, New  
52 York, United States of America;
- 53 <sup>23</sup>Institute of Human Genetics, University of Lübeck, Lübeck, Germany;

54 <sup>24</sup>Department of Human Genetics, University Hospital of Heidelberg, Heidelberg,

55 Germany;

56 + equally contributing first authors

57 \* Corresponding author, [pkrawitz@uni-bonn.de](mailto:pkrawitz@uni-bonn.de)

58

## 59 **Abstract**

60 A large fraction of monogenic disorders causes craniofacial abnormalities with  
61 characteristic facial morphology. These disorders can be diagnosed more efficiently  
62 with the support of computer-aided next-generation phenotyping tools, such as  
63 DeepGestalt. These tools have learned to associate facial phenotypes with the  
64 underlying syndrome through training on thousands of patient photographs. However,  
65 this “supervised” approach means that diagnoses are only possible if the disorder was  
66 part of the training set. To improve recognition of ultra-rare disorders, we created  
67 GestaltMatcher, which uses a deep convolutional neural network based on the  
68 DeepGestalt framework. We used photographs of 17,560 patients with 1,115 rare  
69 disorders to define a “Clinical Face Phenotype Space”. Distance between cases in the  
70 phenotype space defines syndromic similarity, allowing test patients to be matched to  
71 a molecular diagnosis even when the disorder was not included in the training set.  
72 Similarities among patients with previously unknown disease genes can also be  
73 detected. Therefore, in concert with mutation data, GestaltMatcher could accelerate  
74 the clinical diagnosis of patients with ultra-rare disorders and facial dysmorphism, as  
75 well as enable the delineation of novel phenotypes.

## 76 Introduction

77 Rare genetic disorders affect more than 6.2% of the global population<sup>1</sup>. Because  
78 genetic disorders are rare and diverse, accurate clinical diagnosis is a time-consuming  
79 and challenging process, often referred to as the “diagnostic odyssey,<sup>2</sup>” and all  
80 informative clinical features have to be taken into consideration. A large fraction of  
81 patients, particularly those with neurodevelopmental disorders, exhibits craniofacial  
82 abnormalities<sup>3</sup>. If the facial phenotype (“gestalt”) is highly recognizable, such as in  
83 Down syndrome, it may also play an important role in establishing the diagnosis.  
84 Sometimes the gestalt is so characteristic or distinct that it reduces the search space  
85 of candidate genes or can be used to delineate novel phenotype-gene associations<sup>4</sup>.  
86 However, the ability to recognize these syndromic disorders relies heavily on the  
87 clinician’s experience. Reaching a diagnosis is very challenging if the clinician has not  
88 previously seen a patient with an ultra-rare disorder or if the patient presents with a  
89 novel disorder, both of which are increasingly common scenarios.

90 With the rapid development of machine learning and computer vision, a considerable  
91 number of next-generation phenotyping tools have emerged that can analyze facial  
92 dysmorphology using two-dimensional (2D) portraits of patients<sup>5–13</sup>. These tools can  
93 aid in the diagnosis of patients with facial dysmorphism by matching their facial  
94 phenotype with that of known disorders. In 2014, Ferry *et al.* proposed using a Clinical  
95 Face Phenotype Space (CFPS) formed by facial features extracted from images to  
96 perform syndrome classification; the system in that study was trained on photos of  
97 more than 1,500 controls and 1,300 patients with eight different syndromes<sup>5</sup>. Since  
98 then, facial recognition technologies have improved significantly and constitute the  
99 core of the deep-learning revolution in computer vision<sup>14,15</sup>. The current state-of-the-art  
100 framework for syndrome classification, DeepGestalt (Face2Gene, FDNA inc, USA),  
101 has been trained on more than 20,000 patients and currently achieves high accuracy

102 in identifying the correct syndrome for roughly 300 syndromes<sup>12,16</sup>. DeepGestalt has  
103 also demonstrated a strong ability to separate specific syndromes and subtypes,  
104 surpassing human experts' performance<sup>12</sup>. Hence, pediatricians and geneticists  
105 increasingly use such next-generation phenotyping tools for differential diagnostics in  
106 patients with facial dysmorphism. However, most existing tools, including DeepGestalt,  
107 need to be trained on large numbers of photographs, and are therefore limited to  
108 syndromes with at least seven images of different patients. The number of submissions  
109 to diagnostic databases of pathogenic variants, such as ClinVar<sup>17</sup>, has become a good  
110 surrogate for the prevalence of rare disorders. When submissions to ClinVar of disease  
111 genes with pathogenic mutations are plotted in decreasing order, most of the  
112 supported syndromes are on the left, indicating relatively high prevalence (Figure 1).  
113 For instance, Cornelia de Lange syndrome (CdLS), which has been modeled by  
114 multiple tools<sup>5,12</sup>, is caused by mutations in *NIPBL*, *SMC1A*, or *HDAC8*, as well as other  
115 genes, and has been linked to hundreds of reported mutations. However, more than  
116 half of the genes in ClinVar have fewer than ten submissions each (Figure 1). As a  
117 result, most phenotypes have not been modeled because sufficient data are lacking.  
118 Thus, the need to train on large numbers of photographs is a major limitation for the  
119 identification of ultra-rare syndromes.

120 A second limitation of classifiers such as DeepGestalt is that their end-to-end, offline-  
121 trained architecture does not support new syndromes without additional modifications.  
122 In order to model a new syndrome in a deep convolutional neural network (DCNN), the  
123 developer has to go through six separate steps (Supplementary Figure 1), including  
124 collecting images of the new syndrome; changing the classification head, which is the  
125 last layer of the DCNN; retraining the network; and more. In addition, the model cannot  
126 be used to quantify similarities among undiagnosed patients, which is crucial in the  
127 delineation of novel syndromes.

128 A third shortcoming of current approaches is that they are not able to contribute to the  
129 longstanding discussion within the nosology of genetic diseases about  
130 distinguishability. Syndromic differences have been hard to measure objectively<sup>18</sup>, and  
131 decisions to “split” syndromes into separate entities on the basis of perceived  
132 differences or to “lump” syndromes together on the basis of similarities have been  
133 made subjectively. Current tools are unable to quantify the similarities between  
134 syndromes in a way that could shed light on the underlying molecular mechanisms and  
135 guide classification.

136 Our objective is to improve phenotypic decision support for rare disorders. Here we  
137 describe GestaltMatcher, an innovative approach that uses an image encoder to  
138 convert all features of a facial image into a vector of numbers. The encoder can also  
139 be thought of as the penultimate layer of a DCNN that was trained on known  
140 syndromes, such as DeepGestalt. The vectors resulting from the encoder are then  
141 used to build a CFPS for matching a patient’s photo to a gallery of portraits of solved  
142 or unsolved cases. The distance between cases in the CFPS quantifies the similarities  
143 between the faces, thereby matching patients with known syndromes or identifying  
144 similarities between multiple patients with unknown disorders and thereby helping to  
145 define new syndromes. Because GestaltMatcher quantifies similarities between faces  
146 in this way, it addresses all three of the limitations described above: (1) it can identify  
147 “closest matches” among patients with known or unknown disorders, regardless of  
148 prevalence; (2) it does not need new architecture or training to incorporate new  
149 syndromes; and (3) it creates a search space to explore similarity of facial gestalts  
150 based on mutation data, which can point to shared molecular pathways of  
151 phenotypically similar disorders.

## 152 **Results**

153 The feature encoder of GestaltMatcher computes a Facial Phenotypic Descriptor (FPD)  
154 for each portrait image (Figure 2a). Each FPD can be thought of as one coordinate in  
155 the CFPS (Figure 2b). The distances between the FPDs in the CFPS form the basis  
156 for syndrome classification, delineation of novel phenotypes, and patient clustering.  
157 The performance for all three of these use cases depends on the composition of the  
158 training set and the gallery. All benchmarking results described in this section, as well  
159 as those available through the web service, are based on data from Face2Gene (F2G).  
160 The F2G dataset was used to construct a CFPS consisting of 26,065 images from  
161 17,502 subjects who had been diagnosed with a total of 1,115 different syndromes,  
162 each supported by at least two cases. We divided the dataset into two categories, the  
163 *rare* dataset consisting of 816 ultra-rare and novel syndromes, representing  
164 syndromes that we aim to identify, and the *frequent* set, consisting of 299 syndromes  
165 already identified by DeepGestalt. The latter set of known syndromes was also used  
166 to train the encoder. Each category was further split into a gallery (90% of each  
167 syndrome) and a test set (the remaining 10% of each syndrome) (see the Online  
168 methods for details).

169 Since F2G data cannot be shared, we compiled the GestaltMatcher database (GMDB),  
170 consisting of 4,306 images from 3,693 subjects with 257 different syndromes. This  
171 second data set is based on 902 publications, and further cases for which we obtained  
172 consent for sharing. All findings described in this section that are based on the F2G  
173 data can be reproduced qualitatively on the GMDB data and are listed in the  
174 Supplemental Material.

### 175 **Training on images of dysmorphism improves the performance of the FPD**

176 To investigate the importance of using a syndromic features encoder rather than a  
177 normal facial features encoder, we compared FPDs that are based on the same  
178 architecture, but trained on different data. The first encoder, which we refer to as *Enc-*



179 *healthy*, was only trained on data from healthy individuals in CASIA-WebFace<sup>19</sup>. The  
180 second encoder, which we refer to as *Enc-F2G* (for Face2Gene), was first trained on  
181 the faces of healthy subjects and then fine-tuned by training on dysmorphic faces from  
182 the gallery of patients with frequent syndromes. All images were encoded separately  
183 for each encoder. We then evaluated the performance of the encoders on test sets of  
184 syndromes from the frequent set and from the rare set. The performance metric was  
185 the percentage of test cases (with known diagnosis) for which an FPD with the  
186 matching disorder was within the  $k$  closest diagnoses in the CFPS (the top- $k$  accuracy).  
187 The features created by Enc-F2G performed better in the matching process than those  
188 created with Enc-healthy (Table 1). This emphasizes the importance of training the  
189 encoder on data from faces with dysmorphic phenotypes and not only on healthy faces.  
190 The features created by Enc-F2G improved the accuracy of matching within the top-  
191 10 closest images from 31.46% to 49.12% for the frequent category. Furthermore, the  
192 top-10 accuracy improved from 21.77% to 29.56% for the rare syndromes, which do  
193 not overlap with the frequent syndromes. The larger relative improvement of 56% on  
194 the frequent test set versus 36% for the rare set could possibly be explained as Enc-  
195 F2G being better suited to encode syndromes of the frequent set because it was  
196 previously trained on these disorders. Likewise, for some of the 816 novel disorders,  
197 the characteristic features were not yet optimally represented by Enc-F2G because  
198 features of these disorders were not part of the training set.

199 The same trend of improvement by fine-tuning on a diverse but smaller set of  
200 syndromic photos is also seen on the public GMDB dataset (Enc-GMDB vs Enc-F2G  
201 in Supplementary Table 1). These results suggest that an encoder that is fine-tuned  
202 on as many syndromic faces as possible, such as DeepGestalt, is a better fit for the  
203 task of syndrome classification than one trained only on healthy faces. Moreover,  
204 DeepGestalt's FPD provides a better generalization or clustering than the FPD  
205 encoded by CASIA for rare syndromes that it had not previously seen.

## 206 **Syndromic diversity improves the performance on novel disorders**

207 Earlier definitions of the FPD were mainly based on training a network with a small  
208 selection of common and highly characteristic syndromes<sup>5,9</sup>. In principle, we could train  
209 GestaltMatcher's encoder on all 1,115 different syndromes in our dataset. However,  
210 most of the facial phenotypes that have recently been linked to a gene are either ultra-  
211 rare or less distinctive, and using a very unbalanced training set with many ultra-rare  
212 disorders linked to only few cases may add noise without substantial additional benefit.  
213 We therefore analyzed the influence of the number of syndromes on the encoder's  
214 fine-tuning by incrementally increasing their number starting with the most frequent  
215 ones. Due to the imbalance among the disorders added each time, the improvement  
216 could be affected by the additional number of training subjects. Therefore, we used the  
217 same number of subjects for each syndrome. In this section, the test set consists only  
218 of disorders from the rare set that the encoder has not seen. The training procedure  
219 and averaging of the readout is described in detail in the Online methods.

220 When we increase the number of training syndromes, the accuracy increases (Figure  
221 3). In general, the performance is also higher when more individuals per syndrome are  
222 used for training. Particularly when more than 50 syndromes are used, the curve for  
223 training with 20 subjects/syndrome is above the curve for 10 subjects/syndrome, and  
224 so on. The same trend is also shown in the public GMDB dataset (Supplementary  
225 Figures 2 and 3).

226 Moreover, double the number of syndromes is better than double the number of  
227 subjects in most of the combinations (Supplementary Figure 4). The effect of doubling  
228 the number of syndromes used for training is greater when the base sample size is  
229 larger than 1200 subjects (Supplementary Figures 5 and 6). Therefore, both of the  
230 findings suggest that increasing the syndromic diversity in the training set improves the  
231 performance on novel disorders.

## 232 **Top-10 accuracy plateaus when encoders are fine-tuned on more than 150** 233 **syndromes**

234 In the previous section, we analyzed the impact of syndromic diversity in a balanced  
235 setting, that is, the dynamics of increasing the number of syndromes while keeping the  
236 size of the increments (the number of added subjects) equal. In this section we analyze  
237 the influence of the number of syndromes on model training in the real-world scenario;  
238 that is, when using all of the subjects per syndrome (Supplementary Figure 7). The  
239 top-10 accuracy improved considerably until about 150 syndromes, representing  
240 roughly 90% of the subjects in the entire training set. Almost doubling the number of  
241 syndromes to 299 with the remaining 10% of subjects only increases the performance  
242 marginally. From these dynamics, we can conclude that including additional  
243 syndromes beyond 299 for defining the FPD will provide little benefit, and we decided  
244 to proceed with the Enc-F2G encoder in the following section that is based on the 299  
245 syndromes described in the original DeepGestalt paper.

## 246 **Performance comparison between GestaltMatcher and DeepGestalt**

247 To validate the GestaltMatcher approach, we first worked with the 323 images of  
248 patients with 91 syndromes from the London Medical Database (LMD)<sup>20</sup> that were  
249 already used for benchmarking the performance of DeepGestalt<sup>12</sup>. When using the  
250 frequent gallery, which contains syndromes that DeepGestalt currently supports,  
251 GestaltMatcher achieved 64.30% and 86.59% accuracy within the top-10 and top-30  
252 ranks, respectively, which was lower than the 81.28% top-10 accuracy and 88.34%  
253 top-30 accuracy achieved by DeepGestalt with a Enc-F2G softmax approach  
254 (Supplementary Table 2 and 3). However, when we used the gallery of all 1,115  
255 syndromes for GestaltMatcher (frequent + rare) which is a search space that is roughly  
256 four times larger, the top-10 and top-30 dropped by only 2.40 percentage points and  
257 5.17 percentage points, respectively (Supplementary Table 2). Moreover, we

258 performed the same evaluation on the F2G-frequent test set and GMDB-frequent test  
259 set. When the number of syndromes in the gallery was increased from 299 to 1,115,  
260 the top-10 and top-30 also dropped slightly by 2.27 and 3.77 percentage points for the  
261 F2G-frequent test set (Table 1). The results of the GMDB frequent test also dropped  
262 slightly while supporting more than twice the number of syndromes (Supplementary  
263 Table 1). These results indicate that the GestaltMatcher clustering approach is highly  
264 scalable and robust to adding new disorders, without the limitations of a classification  
265 approach.

### 266 **Matching undiagnosed patients from unrelated families**

267 In the second use case, we envision GestaltMatcher as a phenotypic complement to  
268 GeneMatcher<sup>21</sup>. To prove that we can match patients from unrelated families who have  
269 the same disease by using only their facial photos, we selected syndromes from 15  
270 recent GeneMatcher publications with titles containing the phrase “facial  
271 dysmorphism”. In contrast to the benchmarking of the previous section, the gallery now  
272 consists of subjects with rare syndromes to simulate undiagnosed subjects and as a  
273 consequence, ranks refer to individuals and not disorders. For the evaluation we still  
274 have to reveal in the end whether an individual from the gallery is a match for a test  
275 case or not. This implies that non-matching cases can harm the performance more  
276 than in the previous section. For instance, if the first matching individual is at rank 30,  
277 but the 29 non-matching individuals with higher similarity to the test case all together  
278 have only four non-matching disorders, then this match would contribute to the top-5  
279 accuracy in the previous section that matched on disorders but to the top-30 accuracy  
280 in this section that matches to individuals. Only the top-1 accuracy remains the same  
281 in both benchmarks.

282 In this scenario, we matched 30 of 91 subjects and connected 26 of 79 families when  
283 using the top-10 criterion (Table 2 and Supplementary Figure 8). When using the top-

284 30 rank, 48 of 91 subjects were matched, and 40 of 79 families were connected. Enc-  
285 healthy, which is trained only with healthy subjects, matched only 40 out of 91 subjects  
286 and connected 34 out of 79 families using the top-30 rank (Supplementary Table 4).  
287 Hence, using the encoder trained with facial dysmorphic subjects improves the  
288 matching considerably.

289 As an example, in a study of *TMEM94*<sup>22</sup>, eight of the ten photos in six different families  
290 were matched, and five of six families were connected within the top-10 rank. When  
291 the three test images in family 2 (F-2-5, F-2-7, F-2-9) were tested, the other five families  
292 were among those in the top-30 rank (Figure 4). The youngest brother, F-2-5, matched  
293 families 1, 3, 5, and 6, and one sister, F-2-7, matched families 1, 4, and 6. Another  
294 sister, F-2-9, matched families 1, 4, 5, and 6. The six families were recruited at five  
295 different institutes in India, Qatar, the United States (NIH Undiagnosed Diseases  
296 Network), and Switzerland, indicating that GestaltMatcher can also connect patients of  
297 different ethnic origins. However, a more systematic analysis of pairwise distances still  
298 revealed considerably smaller distances between subjects with *de novo* mutations and  
299 their family members than between these subjects and unrelated individuals  
300 (Supplementary Figure 9). This reflects similarities in the nonclinical features of the  
301 face, which is also higher within the same ethnicity and is a known confounding factor  
302 for the GestaltMatcher approach. However, it is a bias that can be attenuated<sup>23</sup> and  
303 will also diminish over time when more diverse training data become available<sup>24</sup>.

#### 304 **GestaltMatcher and human experts agree on syndrome distinctiveness**

305 We hypothesized that some of the ultra-rare disorders that were linked to their disease-  
306 causing genes early on, such as Schuurs-Hoeijmakers syndrome in 2012,<sup>25</sup> have  
307 particularly distinctive facial phenotypes. To systematically analyze the dependence of  
308 disease-gene discovery on the distinctiveness of a facial gestalt, we asked three expert  
309 dysmorphologists (S.M., N.E., and K.W.G.) to grade 299 syndromes on a scale from 1

310 to 3. The more easily they could distinguish the diseases, and the more characteristic  
311 of the disease they deemed the facial features, the higher the score. All three  
312 syndromologists agreed on the same score for 195/299 syndromes, yielding a  
313 concordance of 65.2%. We then selected 50 syndromes as a test set and trained the  
314 model with the remaining 249 syndromes. We analyzed the correlation of the mean of  
315 the distinctiveness score from human experts with the top-10 accuracy that  
316 GestaltMatcher achieves for these syndromes without having been trained on them  
317 (Figure 5a, Supplementary Table 6). The Spearman's rank correlation coefficient was  
318 0.400 ( $P = 0.004$ ), indicating a clear positive correlation between distinctiveness score  
319 and top-10 accuracy. Syndromes with a higher average score tended to perform better,  
320 with Schuurs-Hoeijmakers syndrome being amongst the best-performing syndromes  
321 in GestaltMatcher. The analysis on 20 selected syndromes from the GMDB dataset  
322 also showed a positive correlation between distinctiveness score and top-5 accuracy  
323 (Supplementary Figure 10 and Supplementary Table 7).

324 The correlation for GestaltMatcher accuracy and disease prevalence was not  
325 significant ( $P = 0.130$ ; Figure 5b). This also means that ultra-rare disorders share a  
326 similar distribution of distinctiveness with more common ones, which is important for  
327 estimates about the performance of GestaltMatcher on novel phenotypes in the real  
328 world.

### 329 **Characterization of phenotypes in the CFPS**

330 When syndromologists cannot find a molecular cause for a patient's phenotype in  
331 diagnostic-grade genes after extensive work up in the lab, it becomes a research case  
332 and they may compare the patient's condition to known disorders. For example a  
333 potentially novel phenotype could be described as "syndrome XY-like" to build a case  
334 group for further molecular analysis through genome sequencing. In GestaltMatcher,

335 this is the third use case, and such comparisons can be supported by cluster analysis  
336 in the CFPS with the cosine distance as a similarity metric (Supplementary Table 8).

337 If a novel disease gene has been identified and the similarities of the patients to known  
338 phenotypes outweigh the differences, OMIM groups them into a phenotypic series. On  
339 the gene or protein level, such phenotypic series often correspond to molecular-  
340 pathway diseases, such as GPI-anchor deficiencies for hyperphosphatasia with mental  
341 retardation syndrome (HPMRS) or cohesinopathies for CdLS. For our cluster analysis,  
342 we sampled subjects in our database with subtypes of four large phenotypic series and  
343 found high intersyndrome separability in addition to considerable intrasyndrome  
344 substructure in Noonan syndrome, CdLS, Kabuki syndrome, and  
345 mucopolysaccharidosis. A  $t$ -SNE<sup>26</sup> projection of the FPDs into two dimensions yielded  
346 the best visualization results (Supplementary Figure 11). Although any projection into  
347 a smaller dimensionality might cause a loss of information, the clusters are still clearly  
348 visible for the 743 subjects sampled from these four phenotypic series. This  
349 observation provides further evidence that characteristic phenotypic features are  
350 encoded in the FPDs.

351 To demonstrate the separability of syndromes with facial dysmorphism, we also used  
352  $t$ -SNE to project 4,353 images of the ten syndromes from the frequent set with the  
353 largest number of subjects and 872 images of ten non-distinct syndromes (syndromes  
354 without facial dysmorphism) into 2D space. In addition, we calculated the Silhouette  
355 index<sup>27</sup> for both of these datasets. The FPDs of the frequent syndromes showed ten  
356 clear clusters of subjects (Supplementary Figure 12), but the  $t$ -SNE projection of  
357 subjects with non-distinct syndromes created no clear clusters. Moreover, the  
358 Silhouette index of the frequent syndromes (0.11) was higher than that of the non-  
359 distinct syndromes (-0.005); the negative Silhouette index indicates poor separation  
360 of the non-distinct syndromes.

## 361 **GestaltMatcher as a tool for clinician scientists**

362 The transition of a research case to a diagnostic case is best described by the process  
363 of matching unrelated patients in the CFPS who share a molecular abnormality until  
364 statistical significance is reached. We illustrate this process for the novel disease gene  
365 *PSMC3* in a demonstration on the GestaltMatcher web service (Supplementary Figure  
366 13, [www.gestaltmatcher.org](http://www.gestaltmatcher.org)). Ebstein *et al.* (not yet published) report 18 patients with  
367 a neurodevelopmental disorder of heterogeneous dysmorphism that is caused by *de*  
368 *novo* missense mutations in *PSMC3*, which encodes a proteasome 26S subunit.  
369 Although not all *PSMC3* patients have the same facial phenotype, the proximity of two  
370 unrelated patients in the CFPS who share the same *de novo* *PSMC3* mutation is  
371 exceptional. Their distance is comparable to the pairwise distances of patients with the  
372 recurring missense mutation R203W in *PACS1*, which is the only known cause of  
373 Schuurs-Hoeijmakers syndrome. On the one hand, the high distinctiveness of these  
374 two *PSMC3* cases with the same mutation allows direct matching by phenotype. On  
375 the other hand, the pairwise similarities of 10 out of 18 patients in the CFPS for which  
376 portraits were available also hints that the protein domains have more than one  
377 function. The previously described scalability of GestaltMatcher makes an exploration  
378 of such similarities in the CFPS possible for any number of cases as soon as they have  
379 been added to the gallery of undiagnosed patients.

## 380 **Discussion**

381 GestaltMatcher's ability to match previously unseen syndromes, that is, those for which  
382 no patient is included in the training set, distinguishes it from other approaches. Since  
383 matching of unseen syndromes is not only of importance for ultra-rare disorders but  
384 can be considered for the discovery of novel diseases, GestaltMatcher could also  
385 speed up the process of delineating new disorders.



386 Importantly, GestaltMatcher provides the flexibility to easily scale up the number of  
387 supported syndromes or the number of unsolved cases without substantial loss in  
388 performance. The LMD validation analysis revealed that the use of the softmax  
389 approach, that is classification based on the values of the last layer representing  
390 disorders, outperformed GestaltMatcher. However, the GestaltMatcher encoder, that  
391 is clustering in the CFPS with values of the penultimate layer representing features,  
392 demonstrated high scalability by yielding similar performance when the number of  
393 supported syndromes was increased from 299 to 1,115. Furthermore, the  
394 distinctiveness of a syndrome correlated with the performance (Figure 5a), whereas  
395 syndrome prevalence did not (Figure 5b). Thus, GestaltMatcher can match a syndrome  
396 with a distinguishable facial gestalt even if it is of extremely low prevalence. This  
397 enables us to avoid the long development flow currently required to support and  
398 discover novel syndromes (Supplementary Figure 1). Instead, matching can be offered  
399 instantly for all unsolved cases with available frontal images for which consent has  
400 been provided for inclusion in the tool. If the gallery is populated by cases with a  
401 disease-causing mutation in a diagnostic-grade gene, we consider this a diagnostic  
402 work-up. In contrast, if the gallery is populated by further undiagnosed cases, it is a  
403 use case comparable to GeneMatcher.

404 GestaltMatcher's framework also allows us to abstract the encoding of a dataset away  
405 from the classification task. For example, one can evaluate both phenotypic series and  
406 pleiotropic genes within a single CFPS, or obtain the most-similar patients for each of  
407 the matched syndromes, with minor computational cost (i.e., in real time). Furthermore,  
408 the GestaltMatcher framework computes the similarity between each of the test set  
409 images across the entire dataset of images. This similarity can be computed using  
410 different metrics, e.g., cosine or Euclidean distance. The results are then aggregated  
411 according to the chosen configuration. For example, image similarity can be  
412 aggregated at the patient level or the syndrome level. Furthermore, the dataset can be

413 filtered according to different parameters (such as ethnicity, disease-causing genes, or  
414 age) to further customize the evaluation.

415 One of the key features of GestaltMatcher is the ability to match patients and quantify  
416 their syndromic similarity. For clinician scientists who often face two different tasks in  
417 their daily practice, this means: (1) assessing whether the patient's phenotype is  
418 specific for a known disorder. If e.g. a variant of unclear clinical significance is found in  
419 a diagnostic grade gene, this would be considered as supporting evidence for the  
420 pathogenicity<sup>28,29</sup>. (2) assessing whether the phenotypic similarity of an unsolved case  
421 to other individuals without a diagnosis is high enough to form e.g. a case group that  
422 is further analyzed. This could e.g. result in the identification of potentially deleterious  
423 variants in a novel disease gene and would represent the phenotypic complement to  
424 existing matching approaches on the molecular level. Several online platforms, such  
425 as GeneMatcher, MyGene2 (<https://mygene2.org/MyGene2>), and Matchmaker  
426 Exchange<sup>30</sup>, already allow physicians to look for similar patients based on sequencing  
427 information, and over the past few years these platforms have enabled the matching  
428 of thousands of patients. However, although phenotypic data, encoded e.g. in HPO  
429 terms, are usually exchanged after contact has been established, automated facial  
430 matching technology has not yet been included in any of these platforms.

431 Since its first proof of concept, in which GestaltMatcher was used to identify two  
432 unrelated patients from different countries with the same novel disease, caused by the  
433 same *de novo* mutation in *LEMD2*<sup>4</sup>, our approach has successfully been applied to  
434 further ultra-rare disorders (Figure 1). We matched 40 of 79 different families in 15  
435 GeneMatcher publications by top-30 rank (Table 2 and Supplementary Figure 8), and  
436 11 candidate genes are currently under evaluation. This result shows the power and  
437 potential of GestaltMatcher to identify novel syndromes. Although the number of  
438 individuals and the diversity of their phenotypes will affect the performance, cases with

439 a high syndromic similarity will remain matchable due to the high dimensionality of the  
440 CFPS.

441 We therefore hope that GestaltMatcher will be readily integrated into other matching  
442 platforms to aid in determining which phenotypes should be grouped together into a  
443 syndrome or phenotypic series, as well as linking individual patients to a molecular  
444 diagnosis.

#### 445 **Code availability**

446 GestaltMatcher is a partially proprietary framework. Although the source code for  
447 cropping the face cannot be shared, the architecture of the CNN, as well as a web  
448 service of the trained version of the tool is accessible for use by health care  
449 professionals free of charge at [www.gestaltmatcher.org](http://www.gestaltmatcher.org).

#### 450 **Data availability**

451 The data that support the findings of this study are divided into two groups, sharable  
452 data (GMDB) and non-sharable data (F2G). GMDB is accessible via  
453 [www.gestaltmatcher.org](http://www.gestaltmatcher.org). Restricted data are curated from Face2Gene users under a  
454 license and cannot be published in order to protect patient privacy.

#### 455 **Online methods**

##### 456 **Study approval**

457 This study is governed by the following Institutional Review Board (IRB) approval:  
458 Charité–Universitätsmedizin Berlin, Germany (EA2/190/16); UKB Universitätsklinikum  
459 Bonn, Germany (Lfd.Nr.386/17). The authors have obtained written informed consent  
460 given by the patients or their guardians, including permission to publish photographs.

##### 461 **Face2Gene datasets**

462 We collected images of subjects with clinically or molecularly confirmed diagnoses  
463 from the Face2Gene database (<https://www.face2gene.com>). Extracted, deidentified  
464 data were used to remove poor-quality or duplicated images from the dataset without  
465 viewing the photos. After removing images of insufficient quality, the dataset consisted  
466 of 26,152 images from 17,560 subjects with a total of 1,115 syndromes  
467 (Supplementary Table 9).

468 GestaltMatcher was designed to distinguish syndromes with different properties. We  
469 separated syndromes by the number of affected subjects and whether they had  
470 already been learned by the DeepGestalt model. Supplementary Figure 14 provides  
471 an overview of how the dataset was divided. The current DeepGestalt approach  
472 requires at least seven subjects to learn a novel syndrome. We first used this threshold  
473 to separate the syndromes into “frequent” and “rare” syndromes. The objective of our  
474 study was to improve phenotypic decision support for “rare disorders”. However,  
475 frequent syndromes that are not associated with facial dysmorphic features cannot be  
476 modeled by DeepGestalt. We therefore further selected 299 frequent syndromes that  
477 possess characteristic facial dysmorphism recognized by DeepGestalt as “frequent  
478 syndromes”. The frequent syndromes were used to validate syndrome prediction and  
479 the separability of subtypes of a phenotypic series because these syndromes are  
480 known to have facial dysmorphic features that are well recognized by the DeepGestalt  
481 encoder. For rare syndromes, we sought to demonstrate that GestaltMatcher could  
482 predict a syndrome even if facial images were publicly available for only a few subjects.  
483 It is noteworthy that, for more than half of all known disease-causing genes, fewer than  
484 ten cases with pathogenic variants have been submitted to ClinVar (Figure 1). Of the  
485 1,115 syndromes in the entire dataset, 299 were frequent and 816 were rare.  
486 DeepGestalt cannot yet be applied to rare syndromes.

487 We further divided each of these two datasets into a gallery and a test set. The gallery  
488 is the set of subjects that we intend to match, given a subject from the test set. First,

489 90% of subjects with each frequent syndrome were used to train the models, and the  
490 remaining 10% of subjects were used to validate the DeepGestalt training; the 90%  
491 then became the frequent gallery and the 10% were assigned to the frequent test set.  
492 For the rare dataset, we performed 10-fold cross-validation. In each syndrome, 90%  
493 and 10% of subjects were assigned to the gallery and test set, respectively. The test  
494 sets were designed to have the same distribution of distinctiveness as the training sets.

495 Matching only within a dataset would not represent a real-world scenario. Therefore,  
496 the galleries of the two datasets were later combined into a unified gallery that was  
497 used to search for matched patients.

498 Please note that the threshold of seven subjects to divide the dataset into frequent and  
499 rare is to compare GestaltMatcher to DeepGestalt, which both use the same training  
500 data. We could adjust this threshold higher or even remove this threshold in the future.

#### 501 **GMDB dataset**

502 We collected images of subjects with clinically or molecularly confirmed diagnoses  
503 from publications and individuals that gave appropriate informed consent for the  
504 purpose of this study. This dataset can be used as a public training and test set for  
505 benchmarking and is available at GestaltMatcher Database  
506 (<https://gestaltmatcher.gene-talk.de>).

507 At the time of the data freeze on 9 June 2021, the dataset consisted of 4,306 images  
508 of 3,693 subjects with a total of 257 syndromes from 902 publications (Supplementary  
509 Table 9). Six of the 3,693 subjects have not yet been published, but appropriate  
510 consent has been obtained. For a fair comparison with the Face2Gene dataset, we  
511 performed the data separation in the same way. The dataset was first split by the same  
512 threshold (seven subjects) into frequent and rare datasets, giving 139 syndromes in  
513 the frequent dataset and 118 syndromes in the rare set. Both datasets were also later

514 separated into gallery and test sets. The data split is shown in Supplementary Figure  
515 15. Of the 3,693 subjects in GMDB, 963 are also in Face2Gene dataset. To use the  
516 GMDB rare set as the test set for both the GMDB frequent set and the Face2Gene  
517 frequent set, we made sure that there is no syndrome that is in both the GMDB rare  
518 set and Face2Gene frequent set (Supplementary Figure 16).

### 519 **DeepGestalt encoder**

520 The preprocessing pipeline of DeepGestalt includes point detection, facial alignment  
521 (frontalization), and facial region cropping. During inference, a facial region crop is  
522 forward passed through a deep convolutional network (DCNN) and ultimately gives the  
523 final prediction of the input face image. The DeepGestalt network consists of ten  
524 convolutional layers (Conv) with batch normalization (BN) and a rectified linear  
525 activation unit (ReLU) to embed the input features. After every Conv-BN-ReLU layer,  
526 a max pooling layer is applied to decrease spatial size while increasing the semantic  
527 representation. The classifier part of the network consists of a fully connected linear  
528 layer with dropout (0.5). In this study, we considered the DeepGestalt architecture as  
529 an encoder–classification composition, pipelined during inference. We chose the last  
530 fully connected layer before the softmax classification as the facial feature  
531 representation (facial phenotypic descriptor, FPD), resulting in a vector of size 320.

532 DeepGestalt was first trained on images of healthy individuals from CASIA-WebFace<sup>19</sup>,  
533 and later fine-tuned on a dataset with patient images (Face2Gene or GMDB). The  
534 encoder without fine-tuning on patient images was called Enc-healthy. The encoder  
535 later trained on 299 frequent syndromes in the Face2Gene dataset was named Enc-  
536 F2G. The encoder trained on 139 frequent syndromes in GMDB was named Enc-  
537 GMDB. In the following sections, we have several encoders trained on different  
538 subsets of the Face2Gene and GMDB datasets. The summary of all the encoders used  
539 in this study is shown in Supplementary Table 5. To compare GestaltMatcher and

540 DeepGestalt, we used a model using softmax for predicting syndromes, which we  
541 called “Enc-F2G (softmax)”. This model is the same as Enc-F2G; the only difference  
542 is that Enc-F2G (softmax) used softmax in the last layer for prediction, as in  
543 DeepGestalt, and Enc-F2G used the cosine distance of FPDs for prediction.

544 Our first hypothesis was that images of patients with the same molecularly diagnosed  
545 syndromes or within the same phenotypic series, and who also share similar facial  
546 phenotypes, can be encoded into similar feature vectors under some set of metrics.  
547 Moreover, we hypothesized that DeepGestalt’s specific design choice of using a  
548 predefined, offline-trained, linear classifier could be replaced by other classification  
549 “heads”, for example, *k*-Nearest Neighbors using cosine distance, which we used for  
550 GestaltMatcher.

### 551 **Descriptor projection: Clinical Face Phenotype Space**

552 Each image was encoded by the DeepGestalt encoder, resulting in a 320-dimensional  
553 FPD. These FPDs were further used to form a 320-dimensional space called the  
554 Clinical Face Phenotype Space (CFPS), with each FPD a point located in the CFPS,  
555 as shown in Figure 2. The similarity between two images is quantified by the cosine  
556 distance between them in the CFPS. The smaller the distance, the greater the similarity  
557 between the two images. Therefore, clusters of subjects in the CFPS can represent  
558 patients with the same syndrome, similarities among different disorders, or the  
559 substructure under a phenotypic series.

### 560 **Evaluation**

561 To evaluate GestaltMatcher, we took the images in the test set as input and positioned  
562 them in the CFPS defined by the images of the gallery. We calculated the cosine  
563 distance between each of the test set images (for which the diagnoses were known in  
564 this proof-of-concept study) and all of the gallery images. Then, for each test image, if

565 an image from another subject with the same disorder in the gallery was among the  
566 top- $k$  nearest neighbors, we called it a top- $k$  match. We then benchmarked the  
567 performance by averaging the top- $k$  accuracy (percent of test images with correct  
568 matches within the top  $k$ ) of each syndrome to avoid biasing predictions toward the  
569 major class. We further compared the accuracy of each syndrome in the frequent and  
570 rare syndrome subsets to investigate whether GestaltMatcher can extend DeepGestalt  
571 to support more syndromes. To compare its performance on predicting syndromes with  
572 DeepGestalt, we first performed image aggregation on the syndrome level before  
573 calculating top- $k$  accuracy, which means that only the nearest image of each syndrome  
574 will be taken into account.

#### 575 **London Medical Dataset validation analysis**

576 We compiled 323 images of patients diagnosed with 91 frequent syndromes from the  
577 LMD<sup>19</sup> and used this as the validation set for frequent syndromes. We first evaluated  
578 the validation set using softmax, which is a DeepGestalt method. To compare the  
579 performance with that of GestaltMatcher, we evaluated the performance of  
580 GestaltMatcher on two different galleries: a gallery of frequent syndromes consisting  
581 of 19,950 images of patients with 299 syndromes, and a unified gallery consisting of  
582 22,298 images of patients with 1,115 syndromes. We then reported the top- $k$  accuracy  
583 and compared the results of these three settings (DeepGestalt with softmax,  
584 GestaltMatcher with frequent gallery, and GestaltMatcher with unified gallery).

#### 585 **Rare syndromes analysis**

586 To understand the potential for matching rare syndromes, we trained an encoder,  
587 denoted Enc-F2G-rare, on 467 out of 816 rare syndromes with more than two and  
588 fewer than seven subjects. Ninety percent of the subjects were used to train Enc-F2G-  
589 rare and were later assigned to the gallery. The remaining 10% of subjects were



590 assigned to the test set. We then compared the performance of Enc-F2G-rare and  
591 Enc-F2G using cosine distance and the softmax classifier.

### 592 **Matching undiagnosed patients from unrelated families**

593 We selected 15 articles published from 2015 to 2019 in which GeneMatcher was used  
594 to establish an association of a gene with a novel phenotype with facial dysmorphism  
595 from unrelated families. In total, these studies contained 108 photos of 91 subjects  
596 from 79 families. The details are shown in Table 2. The 15 genes were not among the  
597 Face2Gene frequent syndromes, so we can consider them each as a novel phenotype  
598 to the model. We performed leave-one-out cross-validation on this dataset; that is, we  
599 kept one photo as the test set, and we assigned the rest of the photos to a gallery of  
600 3,533 photos with 816 rare syndromes to simulate the distribution of patients with  
601 unknown diagnosis. We then evaluated the performance by top-1 to top-30 rank. If a  
602 photo of another subject with the same disease-causing gene from an unrelated family  
603 was among the top- $k$  rank, we called it a match.

604 Moreover, we used top- $k$  rank to measure how many unrelated families were  
605 connected. If one unrelated family was among the test photo's top- $k$  rank, the families  
606 were considered to be connected at that rank. How many families were matched to at  
607 least one unrelated family was also represented.

608 When using the GeneMatcher data, we did not perform syndrome aggregation  
609 because aggregation cannot be performed if the syndrome is not known. Instead, we  
610 matched patients rather than predicting disorders.

### 611 **Syndrome facial distinctiveness score**

612 To evaluate the importance of the facial gestalt for clinical diagnosis of the patient, we  
613 asked three dysmorphologists (co-authors Shahida Moosa, Nadja Ehmke, and Karen

614 W. Gripp) to score the usefulness of each syndrome's facial gestalt for establishing a  
615 diagnosis. Three levels were established:

- 616 1. Facial gestalt can be supportive in establishing the clinical diagnosis.
- 617 2. Facial gestalt is important in establishing the clinical diagnosis, but diagnosis  
618 cannot be made without additional clinical features.
- 619 3. Facial gestalt is a cardinal symptom, and a visual or clinical diagnosis is  
620 possible based only on the facial phenotype.

621 We then averaged the grades from the three dysmorphologists for each syndrome.

## 622 **Syndrome prevalence**

623 The prevalence of each syndrome was collected from Orphanet ([www.orpha.net](http://www.orpha.net)). Birth  
624 prevalence was used when the actual prevalence was missing. If only the number of  
625 cases or families was available, we calculated the prevalence by summing the  
626 numbers of all cases or families and dividing by the global population, using 7.8 billion  
627 for the global population and a family size of ten for each family<sup>31</sup>.

## 628 **Unseen syndromes correlation analysis**

629 To investigate the influence of prevalence and distinctiveness score on the  
630 performance of novel syndromes with facial dysmorphism, we selected 50 frequent  
631 syndromes and kept them out of the training set. The 50 syndromes were selected to  
632 have evenly distributed distinctiveness scores and prevalence distribution; the  
633 distributions are shown in Supplementary Figure 17 and Supplementary Table 6. The  
634 encoder (Enc-F2G-exclude-50) was trained on 90% of the subjects from the other 249  
635 frequent syndromes. In addition, we performed random downsampling to remove the  
636 confounding effect of prevalence. For each iteration, we randomly downsampled each  
637 syndrome by assigning five subjects to the gallery and one subject to the test set. We  
638 then averaged the top-10 accuracy of 100 iterations. We calculated Spearman rank

639 correlation coefficients for the following two pairs of data: between top-10 accuracy  
640 and the syndrome's distinctiveness score, and between top-10 accuracy and the  
641 prevalence of syndromes collected from Orphanet.

642 The same analysis was also performed on the GMDB dataset. We selected 20  
643 syndromes from GMDB frequent instead of 50 syndromes because the GMDB dataset  
644 is smaller than the Face2Gene dataset, and we trained the Enc-GMDB-exclude-20 on  
645 the remaining 119 frequent syndromes. The details of the 20 selected syndromes and  
646 the results are reported in Supplementary Table 7. Please note that we report the top-  
647 5 accuracy in the GMDB dataset instead of top-10 accuracy because of the smaller  
648 number of syndromes in the gallery.

#### 649 **Analysis of number of training syndromes and subjects**

650 In this analysis, we evaluated the influence of training with additional syndromes and  
651 subjects to the novel disorders. To avoid an imbalance among the syndromes, we used  
652 the same number of subjects for each syndrome. We first used four different settings  
653 for the number of subjects: 10, 20, 40, and 80. However, not all syndromes have the  
654 four numbers of subjects we mentioned above for training: for 10, 20, 40, and 80  
655 subjects, there are 242, 156, 84, and 40 syndromes. We then defined the ordering of  
656 syndromes we added each time. To add the same syndromes for the four numbers of  
657 subjects each time, we first sorted syndromes with the number of subjects in  
658 descending order. To avoid bias due to having specific disorders added at each position,  
659 we then performed random sorting five times within each of the intervals [1:40], [41,  
660 80], [81, 150], and [151, 240] to generate five different lists of syndromes. Thus, the  
661 ordering from common disorders to rare disorders was by interval rather than by  
662 syndrome. For example, Kabuki syndrome might be in the 9<sup>th</sup> position in the first list,  
663 but in the 20<sup>th</sup> position in the second list, but in each randomly sorted list Kabuki  
664 syndrome is in the first interval.

665 For each of five different lists of training syndromes, we performed the same training  
666 described as follows. We first trained X number of syndromes with ten subjects, where  
667  $X = 10$  to 240, incremented at an interval of ten syndromes. As mentioned above, there  
668 are only 156 syndromes with more than 20 subjects. Thus, we trained syndromes with  
669 20 subjects with  $X = 10$  to 150 syndromes with the same increment of ten syndromes.  
670 We performed the same process for 40 and 80 subjects, with maximums of 80 and 40,  
671 respectively.

672 For each setting (number of subjects, number of syndromes), we had five models. We  
673 then encoded the photos separately with each model and tested them on the rare  
674 syndromes, which had not been seen by the models. In the end, we averaged the  
675 performance by the five models and report the top-10 accuracy for each setting in  
676 Figure 3. We also used the models described above to encode the GMDB dataset,  
677 tested them with the GMDB rare set, and report the results in Supplementary Figure 2.

678 Because the GMDB dataset is smaller than Face2Gene dataset, we were not able to  
679 use the same number of subjects and syndromes to perform the analysis. For the  
680 GMDB dataset, we used 10, 20, 40 for the number of subjects, and the syndrome  
681 intervals of [1, 10], [11, 40], and [41, 80]. The results of training on GMDB and testing  
682 of the GMDB rare set are shown in Supplementary Figure 3.

683 We next wanted to compare two scenarios, double the number of training syndromes  
684 and double the number of training subjects. For example, we first set training on ten  
685 subjects for each of ten syndromes as the base setting, then compared this  
686 performance to training ten subjects for each of 20 syndromes (double syndromes)  
687 and training 20 subjects for each of ten syndromes (double subjects). The base setting  
688 had 100 subjects in total. Double syndromes and double subjects each had 200  
689 subjects. This comparison allows us to understand the different influence of adding

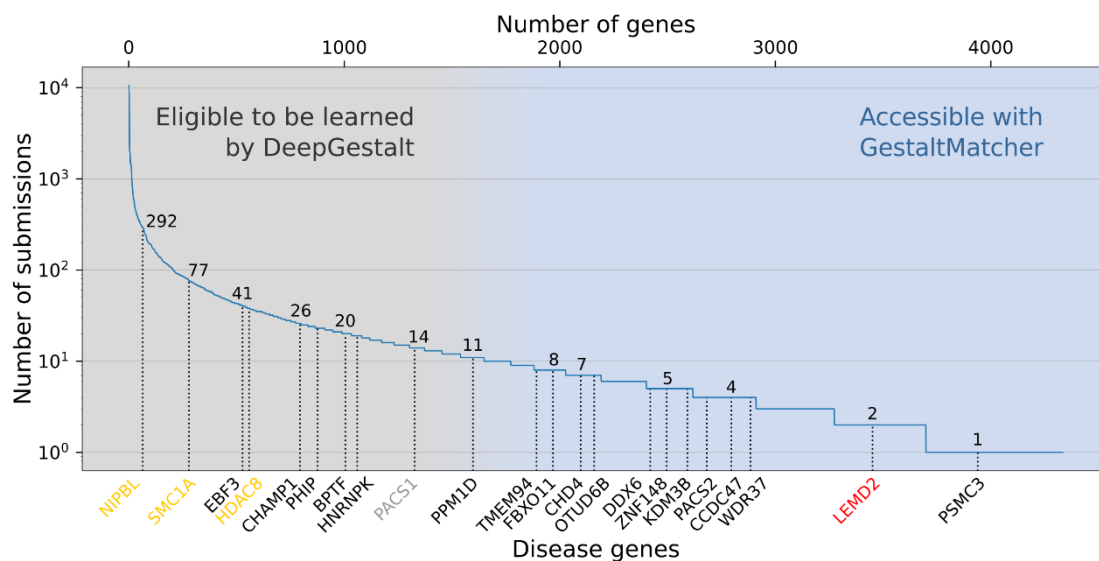
690 more syndromes and adding more subjects. The results are shown in Supplementary  
691 Figures 4-6.

### 692 **Analysis of number of training syndromes in real-world scenario**

693 In this analysis, we trained the encoders with different numbers of syndromes to  
694 simulate the real-world scenario. The difference to the previous section is that we used  
695 all available subjects with each syndrome for the training. To make a fair comparison,  
696 we first used the same ordering of syndromes as in the previous section, and we added  
697 a fifth interval of [241, 299]. For each of the five lists of syndromes, we then trained 16  
698 encoders, each with a different number of training syndromes. The interval of  
699 syndromes was 20 in this analysis due to the long training time. For example, we used  
700 the first ten syndromes in the training list for the first encoder. For the second encoder,  
701 we trained on the first 30 syndromes, and continually increased the number of  
702 syndromes for each subsequent encoder by 20 until we reached 299 syndromes. Thus,  
703 we simulated how syndromes would be included in model training in the real world. We  
704 took the rare syndromes as the test set. We then averaged the performance of five  
705 models with the same number of training syndromes and report the top-10 accuracy in  
706 Supplementary Figure 7.

707

## 708 Figures and tables



709

710 **Figure 1: Subsets of disorders supported by DeepGestalt and GestaltMatcher.**

711 The lower x-axis shows examples of disease genes, and the upper x-axis is the

712 cumulative number of genes. The y-axis shows the number of pathogenic submissions

713 in ClinVar for each gene. The numbers on the curve indicate the number of

714 submissions for each of the indicated genes. Most of the rare disorders that

715 DeepGestalt supports have relatively high prevalence based on their ClinVar

716 submissions, e.g. Cornelia de Lange syndrome (CdLS) which is caused by mutation in

717 *NIPBL*, *SMC1A*, or *HDAC8*, among other genes. Disease genes such as *PACS1* cause

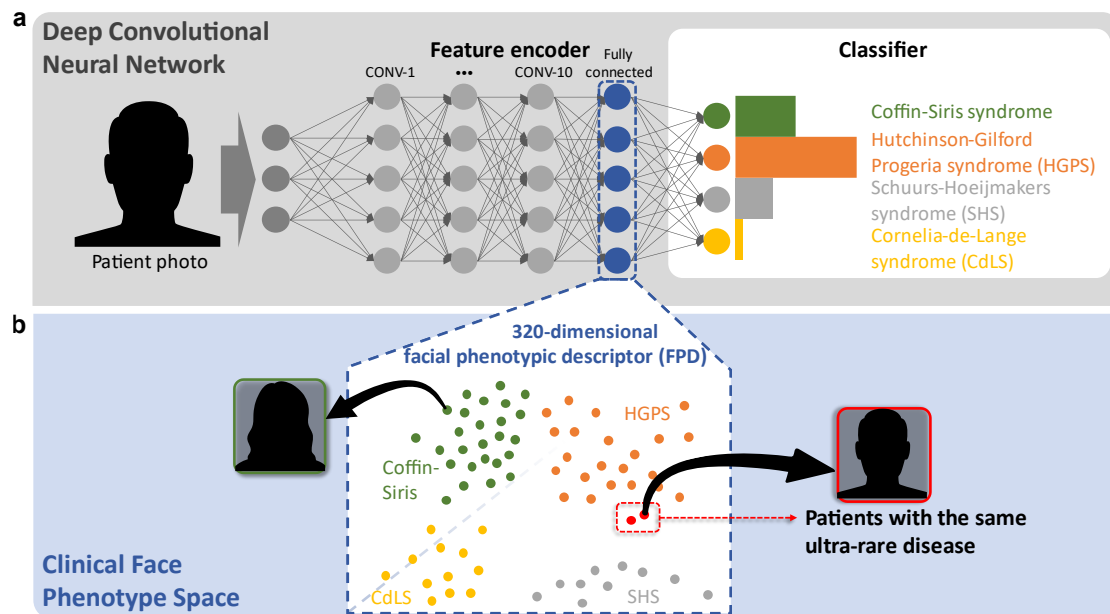
718 highly distinctive phenotypes but are ultra-rare, representing the limit of what current

719 technology can achieve. The first novel disease that was characterized by

720 GestaltMatcher is caused by mutations in *LEMD2*. A candidate disease gene

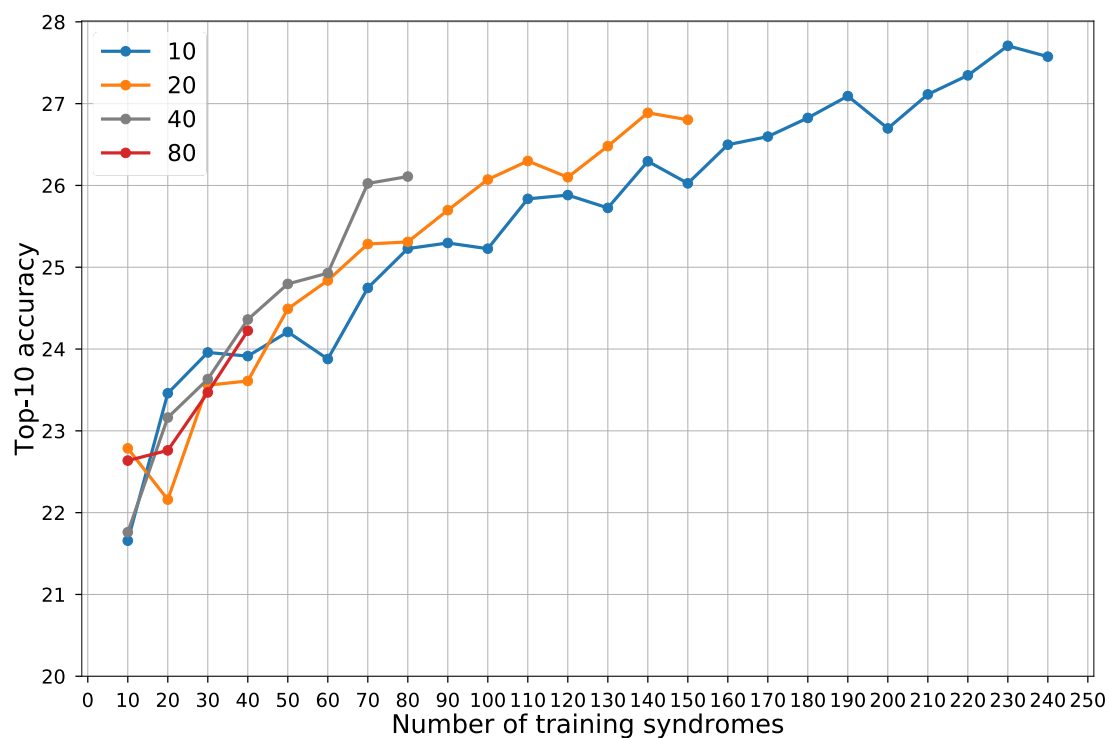
721 associated with a characteristic phenotype that can be identified by GestaltMatcher is

722 *PSMC3*.



723

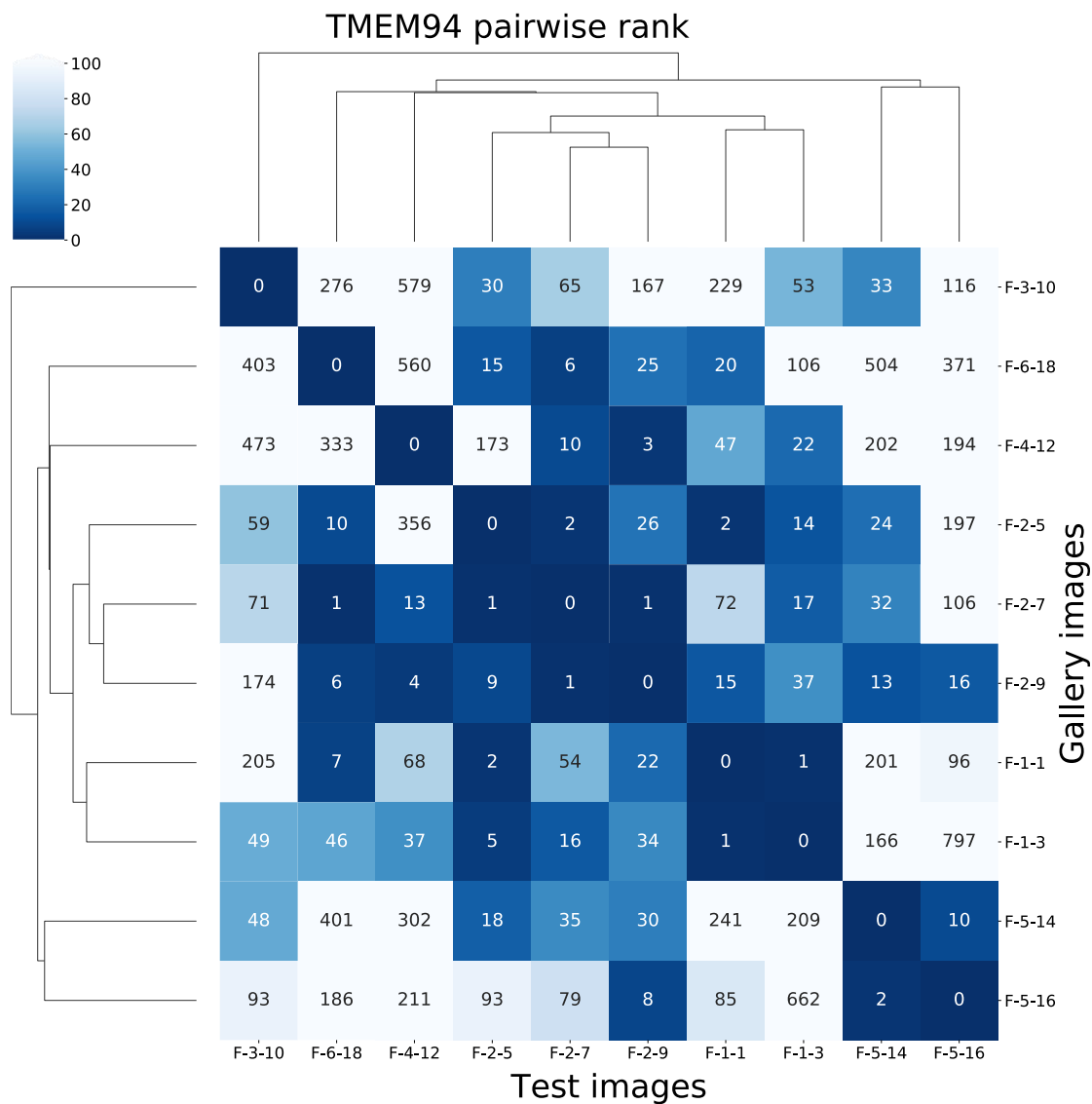
724 **Figure 2: Concept of GestaltMatcher.** **a**, Architecture of a deep convolutional neural  
725 network (DCNN) consisting of an encoder and a classifier. Facial dysmorphic features  
726 of 299 frequent syndromes were used for supervised learning. The last fully connected  
727 layer in the feature encoder was taken as a Facial Phenotypic Descriptor (FPD), which  
728 forms a point in the Clinical Face Phenotype Space (CFPS). **b**, In the CFPS, the  
729 distance between each patient's FPD can be considered as a measure of similarity of  
730 their facial phenotypic features. The distances can be further used for classifying ultra-  
731 rare disorders or matching patients with novel phenotypes. Take the input image as an  
732 example: the patient's ultra-rare disease, which is caused by mutations in *LEMD2*, was  
733 not in the classifier, but was matched with another patient with the same ultra-rare  
734 disorder in the CFPS<sup>4</sup>.



735

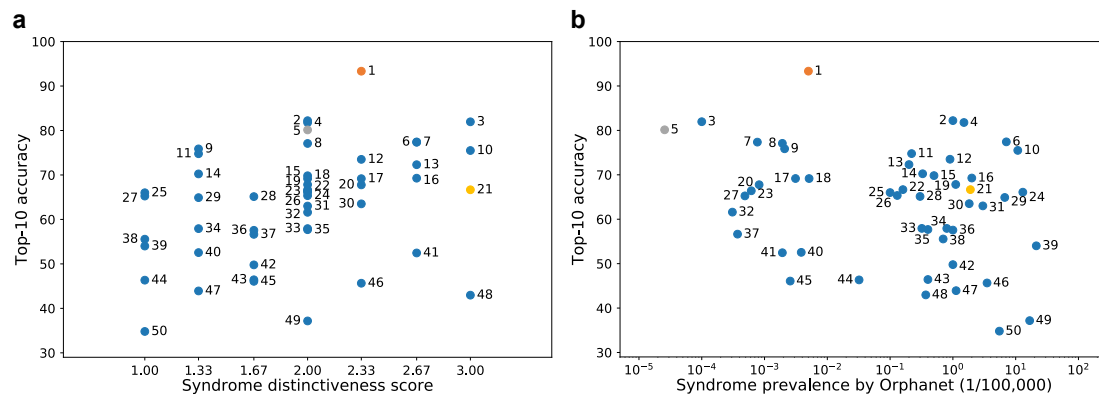
736 **Figure 3: Influence of the number of syndromes included in model training.** The  
737 x-axis is the number of syndromes used in model training. The y-axis shows the  
738 average top-10 accuracy of testing on the rare set. Each line uses the same number  
739 of subjects per syndrome, which is shown in the key. For each point, we train the  
740 models five times with five different splits, and average the results. The null accuracy  
741 (the expected value if the encoder returned random predictions) is 1.2% (10/816).





742

743 **Figure 4: Pairwise ranks of subjects with mutations in *TMEM94*.** Each label  
 744 consists of family numbering and subject numbering, which are the same as in the  
 745 original publication<sup>22</sup>. For example, F-2-7 means the seventh subject in the second  
 746 family. Each column is the result of testing the image indicated at the bottom of the  
 747 column. The number in the box is the rank to the corresponding image in the gallery.  
 748 The fourth column starting from the left is the result of testing F-2-5, and the fourth row  
 749 from the bottom shows that F-1-1 has a rank of 2 for F-2-5. In the fifth to seventh rows  
 750 from the bottom are the ranks from family 2, which is the same family that F-2-5 is from.



751

752 **Figure 5: Correlation among syndrome prevalence, distinctiveness score, and**

753 **top-10 accuracy. a**, Distribution of top-10 accuracy and distinctiveness score. The

754 Spearman rank correlation coefficient was 0.400 ( $P = 0.004$ ). **b**, Distribution of top-10

755 accuracy and prevalence. The Spearman rank correlation coefficient was  $-0.217$  ( $P =$

756 0.130) The details of each syndrome can be found in Supplementary Table 6 using the

757 syndrome ID shown in the figure; syndrome 5 is Schuurs-Hoeijmakers syndrome. The

758 y-axis shows the average top-10 accuracy of the experiments over 100 iterations.

759

760 **Table 1: Performance comparison between classification and clustering with**  
 761 **different encoders on sets of known disorders.**

Test set	Model	Images		Supported syndromes	Null top-1 accuracy	Top-1	Top-5	Top-10	Top-30
		Gallery	Test						
F2G-frequent	Enc-F2G (softmax)	-	2,669	299	0.33%	<b>35.94%</b>	<b>52.45%</b>	<b>63.91%</b>	<b>78.13%</b>
F2G-frequent	Enc-F2G	19,950	2,669	299	0.33%	21.06%	39.62%	49.12%	67.98%
F2G-frequent	Enc-healthy	19,950	2,669	299	0.33%	10.69%	23.69%	31.46%	50.80%
F2G-rare	Enc-F2G	2,348.8	1,183.3	816	0.12%	<b>13.66%</b>	<b>23.62%</b>	<b>29.56%</b>	<b>40.94%</b>
F2G-rare	Enc-healthy	2,348.8	1,183.3	816	0.12%	9.46%	16.87%	21.77%	31.77%
F2G-frequent	Enc-F2G	22,298 <sup>a</sup>	2,669	1,115 <sup>c</sup>	0.09%	<b>20.15%</b>	<b>37.81%</b>	<b>46.85%</b>	<b>64.21%</b>
F2G-frequent	Enc-healthy	22,298 <sup>a</sup>	2,669	1,115 <sup>c</sup>	0.09%	9.70%	22.51%	29.80%	48.24%
F2G-rare	Enc-F2G	22,298.8 <sup>b</sup>	1,183.3	1,115 <sup>c</sup>	0.09%	<b>7.07%</b>	<b>14.19%</b>	<b>17.67%</b>	<b>24.41%</b>
F2G-rare	Enc-healthy	22,298.8 <sup>b</sup>	1,183.3	1,115 <sup>c</sup>	0.09%	4.02%	8.84%	11.73%	16.61%

762 The DCNNs of Enc-F2G (softmax), Enc-F2G, and Enc-healthy have the same architecture.  
 763 Enc-F2G (softmax) and Enc-F2G training were initiated with CASIA-WebFace and further fine-  
 764 tuned on photos of patients in the Face2Gene frequent set. The Enc-F2G (softmax) model is  
 765 the same as Enc-F2G, but using the softmax values of the layer instead of cosine distances  
 766 between the FPDs in the CFPS. For the top-1 to top-30 columns, the best performance in each  
 767 set is boldfaced. The numbers of images and syndromes in the rare set are averaged over ten  
 768 splits. Enc-F2G outperformed Enc-healthy on both types of syndromes, showing the importance  
 769 of fine-tuning on patient photos for learning facial dysmorphic features. The top-10 accuracy of  
 770 Enc-F2G only drops by 2.27 percentage points after increasing the number of cases in the  
 771 gallery and almost quadrupling the number of supported syndromes from 299 to 1,115.

772 <sup>a</sup> Number of images in frequent gallery and rare gallery.

773 <sup>b</sup> Average of ten splits in the frequent gallery and rare gallery.

774 <sup>c</sup> Number of syndromes in the frequent gallery and rare gallery.

775

776 **Table 2: Matching of novel phenotypes on a GeneMatcher validation set.**

Gene	PMID	Total families (Subjects)	Connected families (subjects) <sup>a</sup>	
			Top-10	Top-30
<i>BPTF</i> <sup>32</sup>	28942966	6 (6)	0 (0)	2 (2)
<i>CCDC47</i> <sup>33</sup>	30401460	4 (4)	0 (0)	0 (0)
<i>CHAMP1</i> <sup>34</sup>	27148580	4 (4)	2 (2)	4 (4)
<i>CHD4</i> <sup>35</sup>	27616479	3 (3)	0 (0)	0 (0)
<i>DDX6</i> <sup>36</sup>	31422817	4 (4)	4 (4)	4 (4)
<i>EBF3</i> <sup>37</sup>	28017373	6 (7)	0 (0)	0 (0)
<i>FBXO11</i> <sup>38</sup>	30679813	17 (17)	5 (5)	9 (9)
<i>HNRNPK</i> <sup>39</sup>	26173930	3 (3)	3 (3)	3 (3)
<i>KDM3B</i> <sup>40</sup>	30929739	9 (9)	0 (0)	2 (3)
<i>LEMD2</i> <sup>4</sup>	30905398	2 (2)	2 (2)	2 (2)
<i>OTUD6B</i> <sup>41</sup>	28343629	4 (9)	3 (4)	3 (6)
<i>PACS2</i> <sup>42</sup>	29656858	6 (6)	0 (0)	2 (2)
<i>TMEM94</i> <sup>22</sup>	30526868	6 (10)	5 (8)	6 (10)
<i>WDR37</i> <sup>43</sup>	31327508	4 (4)	2 (2)	3 (3)
<i>ZNF148</i> <sup>44</sup>	27964749	3 (3)	0 (0)	0 (0)
Total	-	79 (91)	26 (30)	40 (48)
Average	-	-	32.91% (32.97%)	50.63% (52.75%)

777 <sup>a</sup> Number of families (subjects) matched by a photo from another family in the top-10 or top-30  
778 rank.

779 In the discovery mode for novel phenotypes, all cases in the gallery are without diagnosis. For  
780 the performance readout, only the correct disease gene of a match is revealed. For individuals  
781 of the *TMEM94* study, e.g. eight out of ten subjects had an image from another family within  
782 the top-10 rank, and five of the six families had at least one subject from another family in their  
783 top-10 rank. For top-30 all subjects and families matched. This table is based on the ranks from  
784 the similarity matrices in Figure 4 and Supplementary Figure 8. The accuracy of connected  
785 subjects corresponds to the accuracy of using Enc-F2G on the F2G-rare test set in the Table 1  
786 in discovery mode in the gallery of almost the same size.

787

## 788      **References**

- 789      1.    Ferreira, C. R. The burden of rare diseases. *Am. J. Med. Genet. A* **179**, 885–892  
790            (2019).
- 791      2.    Baird, P. A., Anderson, T. W., Newcombe, H. B. & Lowry, R. B. Genetic disorders  
792            in children and young adults: A population study. *Am. J. Hum. Genet.* **42**, 677–  
793            693 (1988).
- 794      3.    Hart, T. C. & Hart, P. S. Genetic studies of craniofacial anomalies: clinical  
795            implications and applications. *Orthod. Craniofac. Res.* **12**, 212–220 (2009).
- 796      4.    Marbach, F. *et al.* The Discovery of a LEMD2-Associated Nuclear Envelopathy  
797            with Early Progeroid Appearance Suggests Advanced Applications for AI-Driven  
798            Facial Phenotyping. *Am. J. Hum. Genet.* **104**, 749–757 (2019).
- 799      5.    Ferry, Q. *et al.* Diagnostically relevant facial gestalt information from ordinary  
800            photos. *Elife* **3**, e02020 (2014).
- 801      6.    Kuru, K., Niranjan, M., Tunca, Y., Osvank, E. & Azim, T. Biomedical visual data  
802            analysis to build an intelligent diagnostic decision support system in medical  
803            genetics. *Artif. Intell. Med.* **62**, 105–118 (2014).
- 804      7.    Cerrolaza, J. J. *et al.* Identification of dysmorphic syndromes using landmark-  
805            specific local texture descriptors. in *2016 IEEE 13th International Symposium on*  
806            *Biomedical Imaging (ISBI)* 1080–1083 (2016).
- 807      8.    Wang, K. & Luo, J. Detecting Visually Observable Disease Symptoms from Faces.  
808            *EURASIP J. Bioinform. Syst. Biol.* **2016**, 13 (2016).
- 809      9.    Dudding-Byth, T. *et al.* Computer face-matching technology using two-  
810            dimensional photographs accurately matches the facial gestalt of unrelated  
811            individuals with the same syndromic form of intellectual disability. *BMC Biotechnol.*  
812            **17**, 1–9 (2017).
- 813      10.   Shukla, P., Gupta, T., Saini, A., Singh, P. & Balasubramanian, R. A Deep Learning  
814            Frame-Work for Recognizing Developmental Disorders. in *2017 IEEE Winter*

- 815            *Conference on Applications of Computer Vision (WACV) 705–714 (2017).*
- 816    11. Liehr, T. *et al.* Next generation phenotyping in Emanuel and Pallister-Killian  
817            syndrome using computer-aided facial dysmorphology analysis of 2D photos. *Clin.*  
818            *Genet.* **93**, 378–381 (2018).
- 819    12. Gurovich, Y. *et al.* Identifying facial phenotypes of genetic disorders using deep  
820            learning. *Nature Medicine* vol. 25 60–64 (2019).
- 821    13. van der Donk, R. *et al.* Next-generation phenotyping using computer vision  
822            algorithms in rare genomic neurodevelopmental disorders. *Genet. Med.* **21**,  
823            1719–1725 (2019).
- 824    14. Taigman, Y., Yang, M., Ranzato, M. & Wolf, L. DeepFace: Closing the gap to  
825            human-level performance in face verification. in *Proceedings of the IEEE*  
826            *Computer Society Conference on Computer Vision and Pattern Recognition*  
827            1701–1708 (IEEE Computer Society, 2014).
- 828    15. Huang, G. B., Ramesh, M., Berg, T. & Learned-Miller, E. *Labeled Faces in the*  
829            *Wild: A Database for Studying Face Recognition in Unconstrained Environments.*  
830            <http://vis-www.cs.umass.edu/lfw/>.
- 831    16. Pantel, J. T. *et al.* Advances in computer-assisted syndrome recognition by the  
832            example of inborn errors of metabolism. *J. Inherit. Metab. Dis.* (2018)  
833            doi:10.1007/s10545-018-0174-3.
- 834    17. Landrum, M. J. *et al.* ClinVar: improving access to variant interpretations and  
835            supporting evidence. *Nucleic Acids Res.* **46**, D1062–D1067 (2018).
- 836    18. McKusick, V. A. On lumpers and splitters, or the nosology of genetic disease.  
837            *Perspect. Biol. Med.* **12**, 298–312 (1969).
- 838    19. Yi, D., Lei, Z., Liao, S. & Li, S. Z. Learning Face Representation from Scratch.  
839            (2014).
- 840    20. Winter, R. M. & Baraitser, M. The London Dysmorphology Database. *J. Med.*  
841            *Genet.* **24**, 509–510 (1987).
- 842    21. Sobreira, N., Schiettecatte, F., Valle, D. & Hamosh, A. GeneMatcher: A Matching

- 843 Tool for Connecting Investigators with an Interest in the Same Gene. *Hum. Mutat.*  
844 **36**, 928–930 (2015).
- 845 22. Stephen, J. *et al.* Bi-allelic TMEM94 Truncating Variants Are Associated with  
846 Neurodevelopmental Delay, Congenital Heart Defects, and Distinct Facial  
847 Dysmorphism. *Am. J. Hum. Genet.* **103**, 948–967 (2018).
- 848 23. Alvi, M., Zisserman, A. & Nellåker, C. Turning a blind eye: Explicit removal of  
849 biases and variation from deep neural network embeddings. *Lect. Notes Comput.*  
850 *Sci.* **11129 LNCS**, 556–572 (2019).
- 851 24. Lumaka, A. *et al.* Facial dysmorphism is influenced by ethnic background of the  
852 patient and of the evaluator. *Clin. Genet.* **92**, 166–171 (2017).
- 853 25. Schuurs-Hoeijmakers, J. H. M. *et al.* Recurrent de novo mutations in PACS1  
854 cause defective cranial-neural-crest migration and define a recognizable  
855 intellectual-disability syndrome. *Am. J. Hum. Genet.* **91**, 1122–1127 (2012).
- 856 26. van der Maaten, L. & Hinton, G. Visualizing Data using t-SNE. *J. Mach. Learn.*  
857 *Res.* **9**, 2579–2605 (2008).
- 858 27. Rousseeuw, P. J. Silhouettes: a graphical aid to the interpretation and validation  
859 of cluster analysis. *J. Comput. Appl. Math.* **20**, 53–65 (1987).
- 860 28. Richards, S. *et al.* Standards and guidelines for the interpretation of sequence  
861 variants: A joint consensus recommendation of the American College of Medical  
862 Genetics and Genomics and the Association for Molecular Pathology. *Genet. Med.*  
863 **17**, 405–424 (2015).
- 864 29. Tavtigian, S. V. *et al.* Modeling the ACMG/AMP variant classification guidelines as  
865 a Bayesian classification framework. *Genet. Med.* **20**, 1054–1060 (2018).
- 866 30. Philippakis, A. A. *et al.* The Matchmaker Exchange: A Platform for Rare Disease  
867 Gene Discovery. *Hum. Mutat.* **36**, 915–921 (2015).
- 868 31. Nguengang Wakap, S. *et al.* Estimating cumulative point prevalence of rare  
869 diseases: analysis of the Orphanet database. *Eur. J. Hum. Genet.* **28**, 165–173  
870 (2020).

- 871 32. Stankiewicz, P. *et al.* Haploinsufficiency of the Chromatin Remodeler BPTF  
872 Causes Syndromic Developmental and Speech Delay, Postnatal Microcephaly,  
873 and Dysmorphic Features. *Am. J. Hum. Genet.* **101**, 503–515 (2017).
- 874 33. Morimoto, M. *et al.* Bi-allelic CCDC47 Variants Cause a Disorder Characterized  
875 by Woolly Hair, Liver Dysfunction, Dysmorphic Features, and Global  
876 Developmental Delay. *Am. J. Hum. Genet.* **103**, 794–807 (2018).
- 877 34. Tanaka, A. J. *et al.* De novo pathogenic variants in CHAMP1 are associated with  
878 global developmental delay, intellectual disability, and dysmorphic facial features.  
879 *Cold Spring Harb Mol Case Stud* **2**, a000661 (2016).
- 880 35. Weiss, K. *et al.* De Novo Mutations in CHD4, an ATP-Dependent Chromatin  
881 Remodeler Gene, Cause an Intellectual Disability Syndrome with Distinctive  
882 Dysmorphisms. *Am. J. Hum. Genet.* **99**, 934–941 (2016).
- 883 36. Balak, C. *et al.* Rare De Novo Missense Variants in RNA Helicase DDX6 Cause  
884 Intellectual Disability and Dysmorphic Features and Lead to P-Body Defects and  
885 RNA Dysregulation. *Am. J. Hum. Genet.* **105**, 509–525 (2019).
- 886 37. Harms, F. L. *et al.* Mutations in EBF3 Disturb Transcriptional Profiles and Cause  
887 Intellectual Disability, Ataxia, and Facial Dysmorphism. *Am. J. Hum. Genet.* **100**,  
888 117–127 (2017).
- 889 38. Jansen, S. *et al.* De novo variants in FBXO11 cause a syndromic form of  
890 intellectual disability with behavioral problems and dysmorphisms. *Eur. J. Hum.*  
891 *Genet.* **27**, 738–746 (2019).
- 892 39. Au, P. Y. B. *et al.* GeneMatcher aids in the identification of a new malformation  
893 syndrome with intellectual disability, unique facial dysmorphisms, and skeletal and  
894 connective tissue abnormalities caused by de novo variants in HNRNPK. *Hum.*  
895 *Mutat.* **36**, 1009–1014 (2015).
- 896 40. Diets, I. J. *et al.* De Novo and Inherited Pathogenic Variants in KDM3B Cause  
897 Intellectual Disability, Short Stature, and Facial Dysmorphism. *Am. J. Hum. Genet.*  
898 **104**, 758–766 (2019).



- 899 41. Santiago-Sim, T. *et al.* Biallelic Variants in OTUD6B Cause an Intellectual  
900 Disability Syndrome Associated with Seizures and Dysmorphic Features. *Am. J.*  
901 *Hum. Genet.* **100**, 676–688 (2017).
- 902 42. Olson, H. E. *et al.* A Recurrent De Novo PACS2 Heterozygous Missense Variant  
903 Causes Neonatal-Onset Developmental Epileptic Encephalopathy, Facial  
904 Dysmorphism, and Cerebellar Dysgenesis. *Am. J. Hum. Genet.* **102**, 995–1007  
905 (2018).
- 906 43. Kanca, O. *et al.* De Novo Variants in WDR37 Are Associated with Epilepsy,  
907 Colobomas, Dysmorphism, Developmental Delay, Intellectual Disability, and  
908 Cerebellar Hypoplasia. *Am. J. Hum. Genet.* **105**, 413–424 (2019).
- 909 44. Stevens, S. J. C. *et al.* Truncating de novo mutations in the Krüppel-type zinc-  
910 finger gene ZNF148 in patients with corpus callosum defects, developmental  
911 delay, short stature, and dysmorphisms. *Genome Med.* **8**, 131 (2016).