

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

3 Tzung-Chien Hsieh^{1,+}, Aviram Bar-Haim^{2,+}, Shahida Moosa³, Nadja Ehmke⁴, Karen
4 W. Gripp⁵, Jean Tori Pantel^{1,4}, Magdalena Danyel^{4,6}, Martin Atta Mensah^{4,7}, Denise
5 Horn⁴, Nicole Fleischer², Guilherme Bonini², Alexander Schmid¹, Alexej Knaus¹,
6 Sugirthan Sivalingam¹, Tom Kamphans⁸, Frédéric Ebstein⁹, Elke Krüger⁹, Sébastien
7 Küry^{10,11}, Stéphane Bézieau^{10,11}, Axel Schmidt¹², Sophia Peters¹², Hartmut Engels¹²,
8 Elisabeth Mangold¹², Martina Kreiß¹², Kirsten Cremer¹², Claudia Perne¹², Regina C.
9 Betz¹², Kathrin Grundmann-Hauser¹³, Tobias Haack¹³, Matias Wagner^{14,15}, Theresa
10 Brunet¹⁴, Heidi Beate Bentzen¹⁶, Malte Spielmann¹⁷, Christian Schaaf¹⁸, Stefan
11 Mundlos⁴, Markus M. Nöthen¹², Peter Krawitz^{1,*}

12
13 ¹Institute for Genomic Statistics and Bioinformatics, University Hospital Bonn,
14 Rheinische Friedrich-Wilhelms-Universität Bonn, Bonn, Germany;

15 ²FDNA Inc., Boston, MA, United States;

16 ³Division of Molecular Biology and Human Genetics, Stellenbosch University and
17 Medical Genetics, Tygerberg Hospital, Tygerberg, South Africa;

18 ⁴Institute of Medical Genetics and Human Genetics, Charité-Universitätsmedizin Berlin,
19 Humboldt-Universität zu Berlin and Berlin Institute of Health, Berlin, Germany;

20 ⁵A.I. DuPont Hospital for Children/Nemours, Wilmington, DE, USA;

21 ⁶Berlin Center for Rare Diseases, Charité-Universitätsmedizin Berlin, Humboldt-
22 Universität zu Berlin and Berlin Institute of Health, Berlin, Germany;

23 ⁷Berlin Institute of Health (BIH), Berlin, Germany;

24 ⁸GeneTalk, Bonn, Germany;

25 ⁹Institut für Medizinische Biochemie und Molekularbiologie (IMBM),
26 ¹⁰Universitätsmedizin Greifswald, Greifswald, 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NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.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 ¹⁰CHU Nantes, Service de Génétique Médicale, Nantes, France;

28 ¹¹l'Institut du Thorax, INSERM, CNRS, Université de Nantes, Nantes, France;

29 ¹²Institute of Human Genetics, University of Bonn, Medical Faculty & University
30 Hospital Bonn, Bonn, Germany;

31 ¹³Institute of Medical Genetics and Applied Genomics, University of Tübingen,
32 Tübingen, Germany;

33 ¹⁴Institute of Human Genetics, School of Medicine, Technical University Munich,
34 Munich, Germany;

35 ¹⁵Institute of Neurogenomics, Helmholtz Zentrum München GmbH, German Research
36 Center for Environmental Health, Neuherberg, Germany;

37 ¹⁶Norwegian Research Center for Computers and Law, Faculty of Law, University of
38 Oslo, Oslo, Norway;

39 ¹⁷Institute of Human Genetics, University of Lübeck, Lübeck, Germany;

40 ¹⁸Department of Human Genetics, University Hospital of Heidelberg, Heidelberg,
41 Germany;

42 + equally contributing first authors

43 * Corresponding author, pkrawitz@uni-bonn.de

44 **Abstract**

45 The majority of monogenic disorders cause craniofacial abnormalities with
46 characteristic facial morphology. These disorders can be diagnosed more efficiently
47 with the support of computer-aided next-generation phenotyping tools, such as
48 DeepGestalt. These tools have learned to associate facial phenotypes with the
49 underlying syndrome through training on thousands of patient photographs. However,
50 this “supervised” approach means that diagnoses are only possible if they were part of
51 the training set. To improve recognition of ultra-rare diseases, we created
52 GestaltMatcher, which uses a deep convolutional neural network based on the
53 DeepGestalt framework. We used photographs of 21,836 patients with 1,362 rare

54 disorders to define a “Clinical Face Phenotype Space”. Distance between cases in the
55 phenotype space defines syndromic similarity, allowing test patients to be matched to
56 a molecular diagnosis even when the disorder was not included in the training set.
57 Similarities among patients with previously unknown disease genes can also be
58 detected. Therefore, in concert with mutation data, GestaltMatcher could accelerate
59 the clinical diagnosis of patients with ultra-rare disorders and facial dysmorphism.

60 Introduction

61 Rare genetic disorders affect more than 6.2% of the global population¹. Because
62 genetic disorders are rare and diverse, accurate clinical diagnosis is a time-consuming
63 and challenging process, often referred to as the “diagnostic odyssey.”² Craniofacial
64 abnormalities are present in 30–40% of genetic disorders³. Patients with these
65 syndromic disorders usually have recognizable facies, such as the typical features
66 associated with Down syndrome or Fragile X syndrome. Hence, the facial
67 manifestation can provide a crucial visual hint to help a clinician identify possible
68 underlying disorders, which reduces the search space of candidate genes and speeds
69 up the genetic diagnostic workup. However, the ability to recognize these syndromic
70 disorders relies heavily on the clinician’s experience. Reaching a diagnosis is very
71 challenging if the clinician has not previously seen a patient with an ultra-rare disorder
72 or if the patient presents with a novel disease, both of which are increasingly common
73 scenarios.

74 With the rapid development of machine learning and computer vision, a considerable
75 number of next-generation phenotyping (NGP) tools have emerged that can analyze
76 facial dysmorphology using two-dimensional (2D) portraits of patients^{4–12}. These tools
77 can aid in the diagnosis of patients with facial dysmorphism by matching their facial
78 phenotype with that of known disorders. In 2014, Ferry *et al.* proposed using a Clinical
79 Face Phenotype Space (CFPS) formed by facial features extracted from images to
80 perform syndrome classification; the system in that study was trained on photos of
81 more than 1,500 controls and 1,300 patients with eight different syndromes⁴. Since
82 then, facial recognition technologies have improved significantly and constitute the
83 core of the deep-learning revolution in computer vision^{13,14}. The current state-of-the-art
84 framework for syndrome classification, DeepGestalt, has been trained on more than
85 20,000 patients and currently achieves high accuracy in identifying the correct

86 syndrome for roughly 300 syndromes^{11,15}. DeepGestalt has also demonstrated a
87 strong ability to separate specific syndromes and subtypes, surpassing human experts'
88 performance. Hence, pediatricians and geneticists increasingly use such NGP tools for
89 differential diagnostics in patients with facial dysmorphism. However, most existing
90 tools, including DeepGestalt, need to be trained on large numbers of photographs, and
91 are therefore limited to syndromes with at least seven submissions. The number of
92 submissions to diagnostic databases of pathogenic variants, such as ClinVar¹⁶, has
93 become a good surrogate for the prevalence of rare disorders. When submissions to
94 ClinVar of disease genes with pathogenic mutations are plotted in decreasing order,
95 most of the supported syndromes are on the left, indicating relatively high prevalence
96 (Figure 1). For instance, Cornelia de Lange syndrome (CdLS), which has been
97 modeled by multiple tools^{4,11}, is caused by mutations in *NIPBL*, *SMC1A*, and *HDAC8*,
98 as well as other genes, and has been linked to hundreds of reported mutations.
99 However, more than half of the genes in ClinVar have fewer than ten submissions each
100 (Figure 1). As a result, most phenotypes have not been modeled because sufficient
101 data are lacking. Thus, the need to train on large numbers of photographs is a major
102 limitation for the identification of ultra-rare syndromes.

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

111 A third shortcoming of current approaches is that they are not able to contribute to the
112 longstanding discussion within the nosology of genetic diseases about

113 distinguishability. Syndromic differences have been hard to measure objectively¹⁷, and
114 decisions to “split” syndromes into separate entities on the basis of perceived
115 differences or to “lump” syndromes together on the basis of similarities have been
116 made subjectively. Current tools are unable to quantify the similarities between
117 syndromes in a way that could shed light on the underlying molecular mechanisms and
118 guide classification.

119 Here we describe GestaltMatcher, an innovative approach that uses an image encoder
120 to convert all features of a facial image into a vector of numbers. These vectors are
121 then used to build a CFPS for matching a patient’s photo to a gallery of portraits of
122 solved or unsolved cases. The distance between cases in the CFPS quantifies the
123 similarities between the faces, thereby matching patients with known syndromes or
124 identifying similarities between multiple patients with unknown disorders and thereby
125 helping to define new syndromes. Because GestaltMatcher quantifies similarities
126 between faces in this way, it addresses all three of the limitations described above: (1)
127 it can identify “closest matches” among patients with known or unknown disorders,
128 regardless of prevalence; (2) it does not need new architecture or training to
129 incorporate new syndromes; and (3) it creates a search space to explore similarity of
130 facial gestalts based on mutation data, which can point to shared molecular pathways
131 of phenotypically similar disorders.

132 **Results**

133 The feature encoder of GestaltMatcher computes a Facial Phenotypic Descriptor (FPD)
134 for each portrait image (Figure 2a). Each FPD can be thought of as one coordinate in
135 the CFPS (Figure 2b). The distances between the FPDs in the CFPS form the basis
136 for syndrome classification and patient clustering.

137 The complete dataset used to construct the CFPS consisted of 33,350 images from
138 21,836 subjects who had been diagnosed with a total of 1,362 syndromes, each
139 supported by at least two cases. We divided the dataset into categories of distinct (rare
140 syndromes with facial dysmorphism recognized by DeepGestalt), non-distinct (rare
141 syndromes without described facial dysmorphism, not recognized by DeepGestalt),
142 and target (ultra-rare syndromes with facial dysmorphism that we hope to be able to
143 identify, not evaluable by DeepGestalt). Each category was further split into the gallery
144 (90% of each syndrome) and a test set (the remaining 10% of each syndrome) (see
145 the Online methods for details).

146 **Training on images of dysmorphism improves the performance of the FPD**

147 To investigate the importance of using a syndromic features encoder rather than a
148 normal facial features encoder, we compared FPDs created by the DeepGestalt
149 encoder (Enc-DeepGestalt) with those created by the CASIA-WebFace¹⁸ encoder
150 (Enc-CASIA), which has the same architecture. DeepGestalt was first trained on the
151 faces of healthy subjects and then fine-tuned by training on dysmorphic faces from a
152 gallery of patients with 296 distinct syndromes, whereas Enc-CASIA was trained on
153 the faces of healthy subjects only. All images were encoded separately for each
154 encoder. We then evaluated the performance of the encoders with the distinct, non-
155 distinct, and target test sets. The performance metric was the percentage of test cases
156 (with known diagnosis) for which an FPD with the matching disorder was within the k
157 closest diagnoses in the CFPS (the top- k accuracy). The features created by
158 DeepGestalt performed better in the matching process than those created with Enc-
159 CASIA (Table 1). This emphasizes the importance of training the encoder on data from
160 faces with dysmorphic phenotypes and not only on healthy faces. The features created
161 by DeepGestalt improved the accuracy of matching within the top-10 closest images
162 by 33% for the distinct category. Furthermore, the top-10 accuracy was improved by
163 33% for the target syndromes, which do not overlap with the distinct syndromes. These

164 results suggest that the features encoded by DeepGestalt are a better fit for the task
165 of syndrome classification than the features encoded by the modern CASIA face
166 recognition model. Moreover, DeepGestalt's FPD provides a better generalization than
167 the FPD encoded by CASIA for target syndromes that it had not previously seen.

168 **Top-10 accuracy plateaus when GestaltMatcher is trained on more than 100** 169 **syndromes**

170 Earlier definitions of the FPD were mainly based on training a network with a small
171 selection of common and highly characteristic syndromes^{4,8}. In principle, we could train
172 GestaltMatcher's encoder on all 1,362 different syndromes in our dataset. However,
173 most of the phenotypes that have recently been linked to a gene are either ultra-rare
174 or less distinctive, and using a very unbalanced training set with many ultra-rare
175 disorders linked to only few cases may add noise without substantial additional benefit.
176 We therefore analyzed the influence of the number of syndromes on model training by
177 incrementally increasing their number starting with the most frequent ones (Figure 3).
178 The top-10 accuracy improved with an increase in the number of syndromes until 110
179 syndromes was reached, fluctuated as the number of syndromes further increased to
180 190, and became saturated after 190 syndromes. From these dynamics, we can
181 conclude that including additional syndromes for defining the FPD will provide little
182 benefit, and we decided to model the encoder of GestaltMatcher with the previous 296
183 syndromes of DeepGestalt, rather than all 1,362.

184 **GestaltMatcher performs similar to DeepGestalt with better scalability**

185 To validate the GestaltMatcher approach, we first worked with the 323 images of
186 patients with 90 syndromes from the London Medical Database (LMD)¹⁹ that were
187 already used for benchmarking the performance of DeepGestalt¹¹. When using the
188 distinct gallery, which contains syndromes that DeepGestalt currently supports,
189 GestaltMatcher achieved 74.30% and 89.78% accuracy within the top-10 and top-30

190 ranks, respectively, which was lower than the 84.52% top-10 accuracy and 91.64%
191 top-30 accuracy achieved with DeepGestalt (Table 2 and Supplementary Table 1).
192 However, when we used the gallery of all 1,362 syndromes for GestaltMatcher (distinct,
193 non-distinct, and target), the top-10 and top-30 dropped by only 3.78% and 4.98%,
194 respectively, indicating that the GestaltMatcher approach is highly scalable.

195 **Matching undiagnosed patients from unrelated families**

196 We envision the use of GestaltMatcher as a phenotypic complement to GeneMatcher²⁰.
197 To prove that we can match patients from unrelated families who have the same
198 disease by using only their facial photos, we selected syndromes from 14 recent
199 GeneMatcher publications with a title containing the phrase “facial dysmorphism”. In
200 this test set, we matched 27 of 104 photos and connected 27 of 77 families when using
201 the top-10 criterion (Table 3, Figure 4, and Supplementary Figure 2). When using the
202 top-30 rank, 47 of 104 photos were matched, and 41 of 77 families were connected.
203 Enc-CASIA, which is trained only with healthy subjects, matched only 30 out of 104
204 photos and connected 32 out of 77 families using the top-30 rank (Supplementary
205 Table 2). Hence, using the encoder trained with facial dysmorphic subjects improves
206 the matching considerably.

207 As an example, in a study of *TMEM94*²¹, nine of the ten photos in six different families
208 were matched, and all six families were connected within the top-10 rank. When the
209 three test images in family 2 (F-2-5, F-2-7, F-2-9) were tested, the other five families
210 were among those in the top-30 rank (Figure 4). The youngest family member, F-2-5,
211 matched families 1, 3, 5, and 6, and both the other two family members, F-2-7 and F-
212 2-9, matched families 1, 4, 5, and 6. The six families were recruited at five different
213 institutes in India, Qatar, the United States (NIH Undiagnosed Diseases Network), and
214 Switzerland, indicating that GestaltMatcher can also connect patients of different ethnic
215 origins. However, a more systematic analysis of pairwise distances still revealed

216 considerably smaller distances between subjects with *de novo* mutations and their
217 unaffected family members than between these subjects and unrelated individuals
218 (Supplementary Figure 3). Hence, ethnicity could be a potential confounding factor for
219 the GestaltMatcher approach. However, it is a bias that can be attenuated²² and will
220 also diminish over time when more diverse training data becomes available²³.

221 **Syndrome distinctiveness assessed by GestaltMatcher correlates with expert** 222 **opinion**

223 We hypothesized that ultra-rare disorders that were linked to their disease-causing
224 genes early on, such as Schuurs-Hoeijmakers syndrome in 2012²⁴, have particularly
225 distinctive facial phenotypes. To systematically analyze the dependency of disease-
226 gene discovery on the distinctiveness of a facial gestalt, we asked three expert
227 dysmorphologists to grade 296 syndromes on a scale from 1 to 3. The more easily
228 they could distinguish the diseases, and the more characteristic of the disease they
229 deemed the facial features, the higher the score. All three syndromologists agreed on
230 the same score for 195/296 syndromes, yielding a concordance of 65.8%. We then
231 analyzed the correlation of the mean of this distinctiveness score from human experts
232 with the top-10 accuracy that GestaltMatcher achieves for these syndromes without
233 having been trained on them (Figure 5a). The Spearman's rank correlation coefficient
234 was 0.421 ($P = 0.002$), indicating a clear positive correlation between distinctiveness
235 score and top-10 accuracy. Syndromes with a higher average score tended to perform
236 better, with Schuurs-Hoeijmakers syndrome being amongst the best-performing
237 syndromes in GestaltMatcher. In contrast, there was no significant correlation for
238 GestaltMatcher accuracy and disease prevalence ($P = 0.126$; Figure 5b).

239 **Characterization of phenotypes in the CFPS**

240 When syndromologists cannot reach a final diagnosis for a patient after extensive
241 genetic sequencing, they may compare the patient's condition to a known molecular

242 disorder, for example describing a “syndrome XY-like phenotype”. In GestaltMatcher,
243 such comparisons can be supported by cluster analysis in the CFPS with the cosine
244 distance as a similarity metric (Supplementary Table 3).

245 If a novel disease gene has been identified and the similarities of the patients to known
246 phenotypes outweigh the differences, OMIM groups them into a phenotypic series. On
247 the gene or protein level, such phenotypic series often correspond to molecular-
248 pathway diseases, such as GPI-anchor deficiencies for Hyperphosphatasia with
249 mental retardation syndrome (HPMRS) or cohesinopathies for CdLS. For our cluster
250 analysis, we sampled subjects in our database with subtypes of four large phenotypic
251 series and found high inter-syndrome separability in addition to considerable intra-
252 syndrome substructure in e.g. Noonan syndrome, CdLS, or mucopolysaccharidosis. A
253 *t*-SNE²⁵ projection of the FPDs into two dimensions yielded the best visualization
254 results (Supplementary Figure 4). Although any projection into a smaller dimensionality
255 might cause a loss of information, the clusters are still clearly visible for the 743
256 subjects sampled from these four phenotypic series. This observation provides further
257 evidence that characteristic phenotypic features are encoded in the FPDs.

258 To demonstrate the separability of syndromes with facial dysmorphism, we also used
259 *t*-SNE to project 4,353 images of the ten distinct syndromes with the largest number of
260 subjects and 872 images of ten non-distinct syndromes into 2D space. In addition, we
261 calculated the Silhouette index²⁶ for both of these datasets. The FPDs of the distinct
262 syndromes showed ten clear clusters of subjects (Supplementary Figure 5), but the *t*-
263 SNE projection of subjects with non-distinct syndromes created no clear clusters.
264 Moreover, the Silhouette index of the distinct syndromes (0.11) was higher than that
265 of the non-distinct syndromes (−0.005); the negative Silhouette index indicates poor
266 separation of the non-distinct syndromes.

267 **GestaltMatcher as a tool for clinician scientists**

268 The transition of a research case to a diagnostic case is best described by the process
269 of matching unrelated patients in the CFPS with a shared molecular cause until
270 statistical significance is reached. We illustrate this process for the novel disease gene
271 *PSMC3* in a demonstration on the GestaltMatcher website (Supplementary Figure 6,
272 www.gestaltmatcher.org). Ebstein *et al.* (not yet published) report 18 patients with a
273 neurodevelopmental disorder of heterogeneous dysmorphism that is caused by *de*
274 *novo* missense mutations in *PSMC3*, which encodes a proteasome 26S subunit.
275 Although not all patients have a single facial gestalt in common, the proximity of two
276 unrelated patients in the CFPS who share the same *de novo PSMC3* mutation is
277 exceptional. Their distance is comparable to the pairwise distances of patients with the
278 reoccurring missense mutation R203W in *PACS1*, which is the only known cause of
279 Schuurs-Hoeijmakers syndrome. On the one hand, the high distinctiveness of these
280 two *PSMC3* cases with the same mutation allows direct matching by phenotype. On
281 the other hand, the pairwise similarities of 10 out of 18 patients in the CFPS for which
282 portraits were available, also hints that the protein domains have more than one
283 function. The previously described scalability of GestaltMatcher makes an exploration
284 of such similarities in the CFPS possible for any number of cases as soon as they have
285 been added to the gallery of undiagnosed patients.

286 Discussion

287 GestaltMatcher's ability to match previously unseen syndromes, i.e., those for which
288 no patient is included in the training set, distinguishes it from other tools. Because
289 matching of unseen syndromes can be considered the discovery of novel diseases,
290 GestaltMatcher could speed up the process of defining new diseases.

291 Importantly, GestaltMatcher provides the flexibility to easily scale up the number of
292 supported syndromes. Although the LMD validation analysis revealed that the use of
293 softmax to predict syndromes trained in the model outperformed GestaltMatcher,

294 GestaltMatcher demonstrated high scalability by yielding similar performance when the
295 number of supported syndromes in the CFPS was increased from 296 to 1,362.
296 Furthermore, the distinctiveness of a syndrome correlated with the performance
297 (Figure 5a), whereas syndrome prevalence did not (Figure 5b). Thus, GestaltMatcher
298 can match a syndrome with a distinguishable facial gestalt even if it is of extremely low
299 prevalence. This enables us to avoid the long development flow currently required to
300 support and discover novel syndromes (Supplementary Figure 1). Instead, matching
301 can be offered instantly for all undiagnosed cases with available frontal images for
302 which consent has been provided for inclusion in the tool.

303 GestaltMatcher's framework also allows us to abstract the encoding of a dataset away
304 from the classification task. For example, one can evaluate both phenotypic series and
305 pleiotropic genes within a single CFPS, or obtain the most-similar patients for each of
306 the matched syndromes, with minor computational cost (i.e., in real time). Furthermore,
307 the GestaltMatcher framework computes the similarity between each of the test set
308 images across the entire dataset of images. This similarity can be computed using
309 different metrics, e.g., cosine or Euclidean distance. The results are then aggregated
310 according to the chosen configuration. For example, image similarity can be
311 aggregated at the patient level or the syndrome level. Furthermore, the dataset can be
312 filtered according to different parameters (such as ethnicity, number of affected genes,
313 or age) to further customize the evaluation.

314 One of the most important features of GestaltMatcher is the ability to match patients
315 with highly similar facial features. Clinicians are often faced with the challenge of
316 finding enough patients with a similar phenotype to statistically link the phenotype to a
317 gene. This is especially true when dealing with presumed novel or extremely rare
318 Mendelian disorders. Several online platforms, such as GeneMatcher, MyGene2
319 (<https://mygene2.org/MyGene2>), and Matchmaker Exchange²⁷, allow physicians to
320 look for similar patients based on phenotypic data, such as HPO terms, or genomic

321 sequencing information, and over the past few years, these platforms have facilitated
322 the matching of thousands of patients. However, although facial phenotypes are crucial
323 for allowing physicians to determine whether two patients have a similar disorder,
324 automated facial matching technology has not yet been included in any of these “gene
325 matching” platforms. We expect that GestaltMatcher will be readily integrated into
326 these matching platforms to aid in determining which phenotypes should be grouped
327 together into a syndrome or phenotypic series, as well as linking individual patients to
328 a molecular diagnosis.

329 Since its first proof of concept, in which GestaltMatcher was used to identify two
330 unrelated patients from different countries with the same novel disease, caused by the
331 same *de novo* mutation in *LEMD2*²⁸, our approach has successfully be applied to
332 further ultra-rare disorders (Figure 1). We matched 41 of 77 different families in 14
333 GeneMatcher publications by top-30 rank, and 11 candidate genes are currently under
334 evaluation. This result shows the power and potential of GestaltMatcher to identify
335 novel syndromes.

336 **Online Methods**

337 **Study approval**

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

342 **Datasets**

343 We collected images of subjects with clinically or molecularly confirmed diagnoses
344 from the Face2Gene database (<https://www.face2gene.com>). Extracted, deidentified
345 data were used to remove poor-quality or duplicated images from the dataset without

346 viewing the photos. After removing images of insufficient quality, the dataset consisted
347 of 33,350 images from 21,836 subjects with a total of 1,362 syndromes.

348 GestaltMatcher was designed to distinguish syndromes with different properties. We
349 separated syndromes by the number of affected subjects and whether they had
350 already been learned by the DeepGestalt model. Supplementary Figure 7 provides an
351 overview of how the dataset was divided. The current DeepGestalt approach requires
352 at least seven subjects to learn a novel syndrome. We first used this threshold to
353 separate the syndromes into rare and ultra-rare syndromes. We denoted ultra-rare
354 syndromes as “target” syndromes because the objective of our study was to improve
355 phenotypic decision support for these disorders. However, rare syndromes that are not
356 associated with facial dysmorphic features cannot be modeled by DeepGestalt. We
357 therefore further divided rare syndromes into “distinct” (possessing characteristic facial
358 dysmorphism recognized by DeepGestalt) and “non-distinct” (without facial
359 dysmorphic features or that cannot be recognized by DeepGestalt). The distinct
360 syndromes were used to validate syndrome prediction and the separability of subtypes
361 of a phenotypic series because these syndromes are known to have facial dysmorphic
362 features that are well recognized by the DeepGestalt encoder. We excluded autism
363 from the non-distinct group of syndromes in this study because it had many more
364 subjects than other non-distinct syndromes, leading to an imbalanced dataset. For
365 target syndromes, we sought to demonstrate that GestaltMatcher could predict a
366 syndrome even if facial images were publicly available for only a few subjects. It is
367 noteworthy that, for more than half of all known disease-causing genes, fewer than ten
368 cases with pathogenic variants have been submitted to ClinVar (Figure 1). Of the 1,362
369 syndromes in the entire dataset, 296 were distinct, 242 non-distinct, and 824 target.
370 DeepGestalt cannot yet be applied to non-distinct and target syndromes.

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

373 90% of subjects with each distinct syndrome were used for training models, and the
374 remaining 10% of subjects were used to validate DeepGestalt training; the 90% then
375 became the distinct gallery and the 10% were assigned to the distinct test set. For the
376 target and non-distinct datasets, we performed 10-fold cross-validation. In each
377 syndrome, 90% and 10% of subjects were assigned to the gallery and test set,
378 respectively.

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

382 **DeepGestalt encoder**

383 The preprocessing pipeline of DeepGestalt includes point detection, facial alignment
384 (frontalization), and facial region cropping. During inference, facial region crop is
385 forward passed through a deep convolutional network (DCNN), and ultimately got the
386 final prediction of the input face image. The DeepGestalt network consists of ten
387 convolutional layers (Conv) with batch normalization (BN) and a rectified linear
388 activation unit (ReLU) to embed the input features. After every Conv-BN-ReLU layer,
389 a max pooling layer is applied to decrease spatial size while increasing the semantic
390 representation. The classifier part of the network consists of a fully connected linear
391 layer with dropout (0.5). In this study, we considered the DeepGestalt architecture as
392 an encoder–classification composition, pipelined during inference. We chose the last
393 fully connected layer before the softmax classification as the facial feature
394 representation (facial phenotypic descriptor, FPD), resulting in a vector of size 320.
395 The encoder trained on 296 distinct syndromes was named Enc-DeepGestalt.

396 Our first hypothesis was that images of patients with the same molecularly diagnosed
397 syndromes or within the same phenotypic series, and who also share similar facial
398 phenotypes, can be encoded into similar feature vectors under some set of metrics.

399 Moreover, we hypothesized that DeepGestalt’s specific design choice of using a
400 predefined, offline-trained, linear classifier could be replaced by other classification
401 “heads”, for example, k -Nearest Neighbors using cosine distance, which we used for
402 GestaltMatcher.

403 **Descriptor projection: Clinical Face Phenotype Space**

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

412 **Evaluation**

413 To evaluate GestaltMatcher, we took the images in the test set as input and positioned
414 them in the CFPS defined by the images of the gallery. We calculated the cosine
415 distance between each of the test set images and all of the gallery images. Then, for
416 each test image, if an image from another subject with the same disorder in the gallery
417 was among the top- k nearest neighbors, we called it a top- k match. We then
418 benchmarked the performance by top- k accuracy (percent of test images with correct
419 matches within the top k). We further compared the accuracy of each syndrome in the
420 distinct, non-distinct, and target syndrome subsets to investigate whether
421 GestaltMatcher can extend DeepGestalt to support more syndromes.

422 **London Medical Dataset validation analysis**

423 We compiled 323 images of patients diagnosed with 90 distinct syndromes from the
424 LMD¹⁹ and used this as the validation set for distinct syndromes. We first evaluated
425 the validation set using softmax, which is a DeepGestalt method. To compare the
426 performance with that of GestaltMatcher, we evaluated the performance of
427 GestaltMatcher on two different galleries: a gallery of distinct syndromes consisting of
428 20,091 images of patients with 296 syndromes, and a unified gallery consisting of
429 27,826 images of patients with 1,362 syndromes. We then reported the top-*k* accuracy
430 and compared the results of these three conditions (DeepGestalt with softmax,
431 GestaltMatcher with distinct gallery, and GestaltMatcher with unified gallery).

432 **Target syndromes analysis**

433 To understand the potential for matching target syndromes, we trained an encoder,
434 denoted Enc-Target, on 477 out of 824 target syndromes with more than three and
435 fewer than seven subjects. Ninety percent of the subjects were used to train Enc-
436 Target and were later assigned to the gallery. The remaining 10% of subjects were
437 assigned to the test set. We then compared the performance of Enc-Target and Enc-
438 DeepGestalt (see previous section) using cosine distance and the softmax classifier.

439 **Syndrome facial distinctiveness score**

440 To evaluate the importance of the facial gestalt for clinical diagnosis of the patient, we
441 asked three dysmorphologists to score the usefulness of each syndrome's facial
442 gestalt for establishing a diagnosis. Three levels were established:

- 443 1. Facial gestalt can be supportive in establishing the clinical diagnosis.
- 444 2. Facial gestalt is important in establishing the clinical diagnosis, but diagnosis
445 cannot be made without additional clinical features.
- 446 3. Facial gestalt is a cardinal symptom, and a visual or clinical diagnosis is
447 possible based only on the facial phenotype.

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

449 **Syndrome prevalence**

450 The prevalence of each syndrome was collected from Orphanet (www.orpha.net). Birth
451 prevalence was used when the actual prevalence was missing. If only the number of
452 cases or families was available, we calculated the prevalence by summing the
453 numbers of all cases or families and dividing by the global population, using 7.8 billion
454 for the global population and a family size of ten for each family²⁹.

455 **Unseen syndromes correlation analysis**

456 To investigate the influence of prevalence and distinctiveness score on the
457 performance for novel syndromes with facial dysmorphism, we selected 50 distinct
458 syndromes and kept them out of the training set. The 50 syndromes were selected to
459 have evenly distributed distinctiveness scores and prevalence distribution; the
460 distributions are shown in Supplementary Figure 7 and Table 4. The encoder (Enc-
461 unseen) was trained on 90% of the subjects from the other 246 distinct syndromes. In
462 addition, we performed random downsampling to remove the confounding effect of
463 prevalence. For each iteration, we randomly downsampled each syndrome by
464 assigning five subjects to the gallery and one subject to the test set. We then averaged
465 the top-10 accuracy of 100 iterations. We calculated Spearman rank correlation
466 coefficients for the following two pairs of data: the first between top-10 accuracy and
467 the syndrome's distinctiveness score, and the second between top-10 accuracy and
468 the prevalence of syndromes collected from Orphanet.

469 **Analysis of number of training syndromes**

470 In this analysis, we trained the encoders with different numbers of syndromes. We first
471 sorted the syndromes by the number of subjects in each syndrome, in descending
472 order. We then trained 13 encoders, each with a different number of training syndromes.
473 We used the ten most common syndromes in the training set for the first encoder. For

474 the second encoder, we trained on the top 30 syndromes, and continually increased
475 the number of syndromes for each subsequent encoder by 20 until we reached 246
476 syndromes. Thus, we simulated how syndromes would be included in model training
477 in the real world. We took the 50 selected distinct syndromes as the test set and
478 performed random downsampling as described in the previous section; the only
479 difference was that we used encoders trained from ten to 246 syndromes.

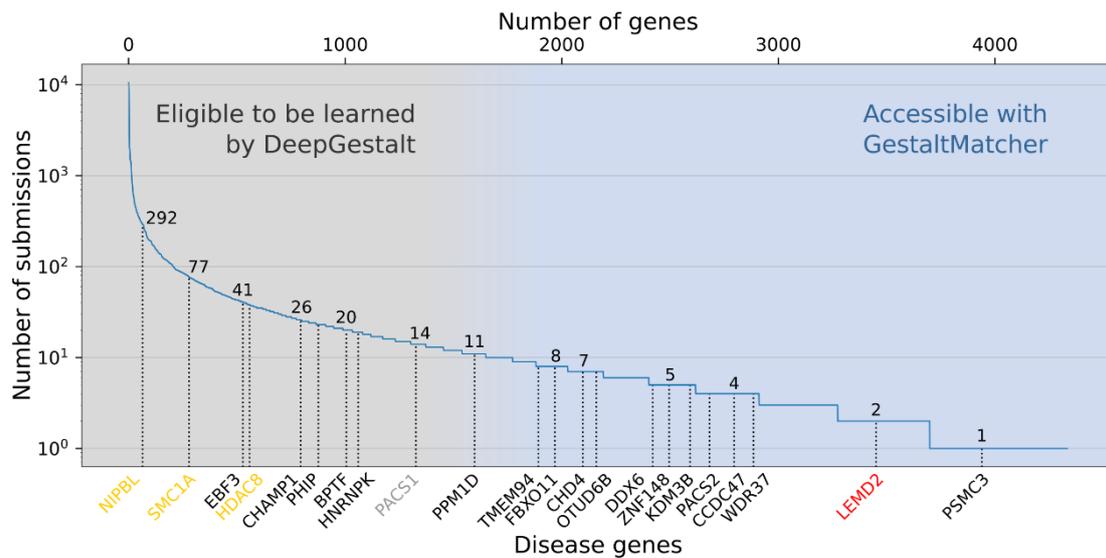
480 **GeneMatcher validation analysis**

481 We selected 14 publications in which GeneMatcher was used to match patients with
482 facial dysmorphism from unrelated families. In total, these studies contained 104
483 photos of 89 subjects from 77 families. The details are shown in Table 3. We performed
484 leave-one-out cross-validation on this dataset, i.e., we kept one photo as the test set,
485 and we assigned the rest of the photos to a gallery of 3,636 photos with 824 target
486 syndromes to simulate the distribution of patients with unknown diagnosis. We then
487 evaluated the performance by top-1 to top-30 rank. If a photo of another subject with
488 the same disease-causing gene from an unrelated family was among the top- k rank,
489 we called it a match.

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

494

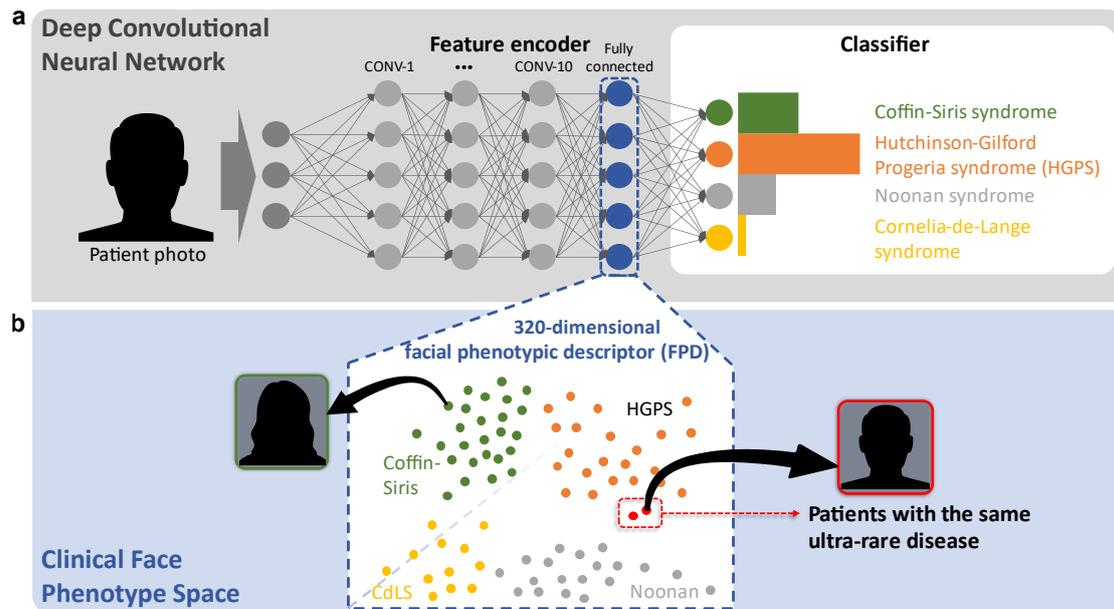
495 **Figures and tables**



496

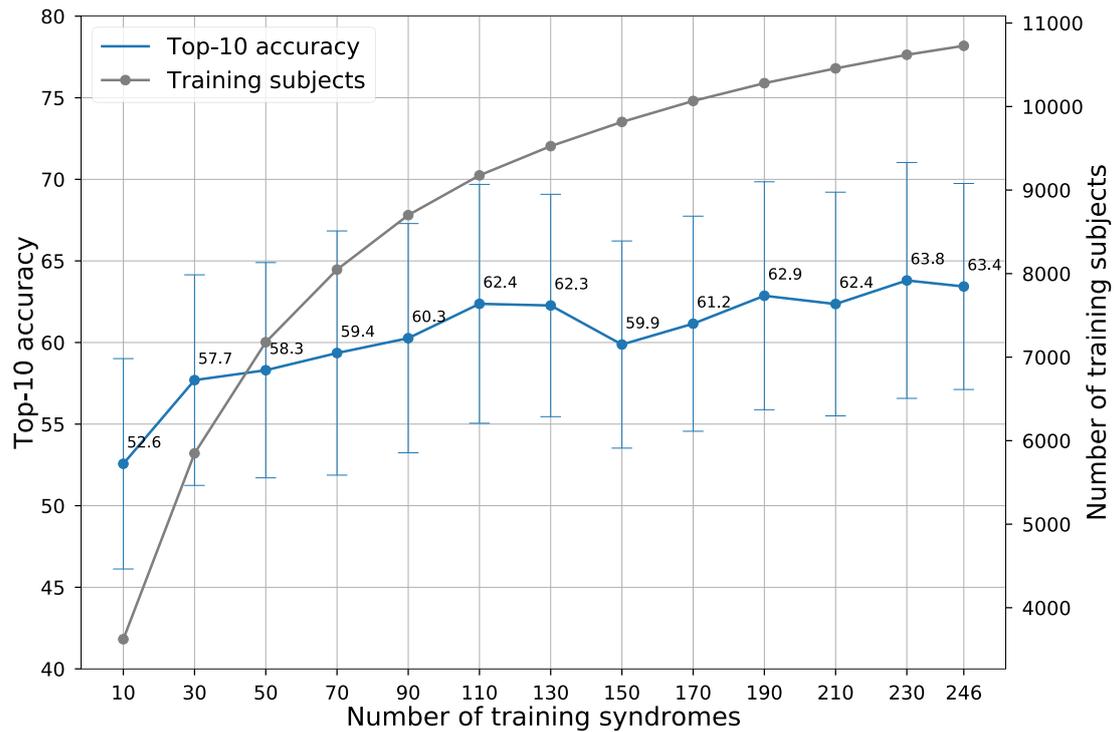
497 **Figure 1: Subsets of disorders supported by DeepGestalt and GestaltMatcher.**

498 The lower x-axis shows examples of disease genes, and the upper x-axis is the
 499 cumulative number of genes. The y-axis shows the number of pathogenic submissions
 500 in ClinVar for each gene. The numbers on the curve indicate the number of
 501 submissions for each of the indicated genes. Most of the rare disorders that
 502 DeepGestalt supports, have a higher prevalence based on their ClinVar submissions,
 503 e.g. Cornelia de Lange syndrome CdLS which is caused by mutation in e.g. *NIPBL*,
 504 *SMC1A*, and *HDAC8*. Disease genes such as *PACS1*, cause highly distinctive
 505 phenotypes but are ultra-rare, representing the limit of what current technology can
 506 achieve. The first novel disease that was characterized by GestaltMatcher, is caused
 507 by mutations in *LEMD2*. A candidate disease gene with a characteristic phenotype
 508 feasible to identification by GestaltMatcher is *PSMC3*.



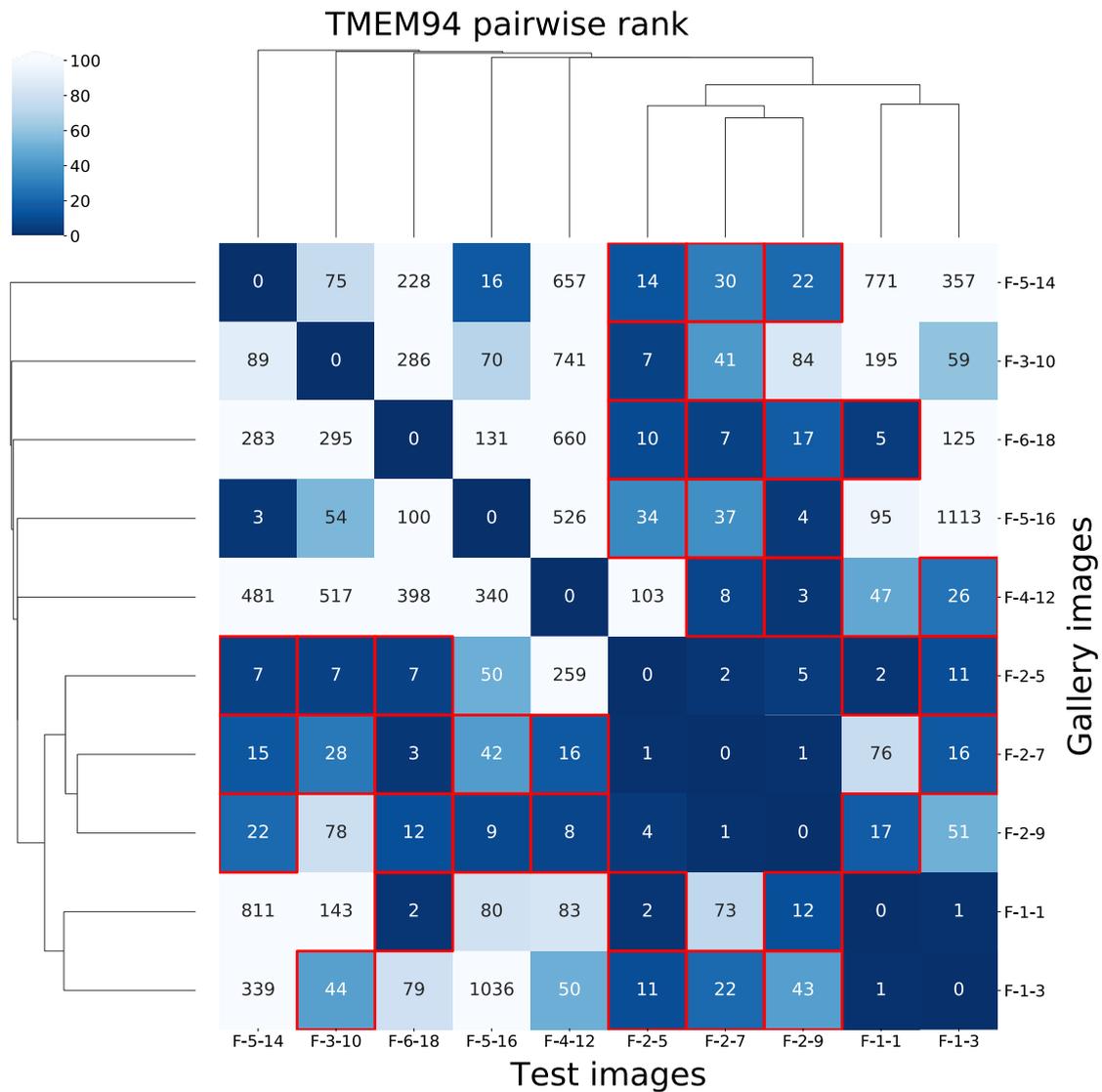
509

510 **Figure 2: Concept of GestaltMatcher.** **a**, Architecture of a deep convolutional neural
511 network (DCNN) consisting of an encoder and a classifier. Facial dysmorphic features
512 of 296 distinct rare syndromes were used for supervised learning. The last fully
513 connected layer in the feature encoder was taken as a Facial Phenotypic Descriptor
514 (FPD), which forms a point in the Clinical Face Phenotype Space (CFPS). **b**, In the
515 CFPS, the distance between each patient's FPD can be considered as a measure of
516 similarity of their facial phenotypic features. The distances can be further used for
517 classifying ultra-rare disorders or matching patients with novel phenotypes. Take the
518 input image as an example: the patient's ultra-rare disease, which is caused by
519 mutations in *LEMD2*, was not in the classifier, but it could match another patient with
520 the same ultra-rare disorder in CFPS.



521

522 **Figure 3: Influence of the number of syndromes included in model training.** The
523 x-axis is the number of syndromes used in model training. The left y-axis shows the
524 average top-10 accuracy over 100 iterations, and the error bars show standard
525 deviation. The right y-axis is the cumulative number of subjects in the training
526 syndromes.

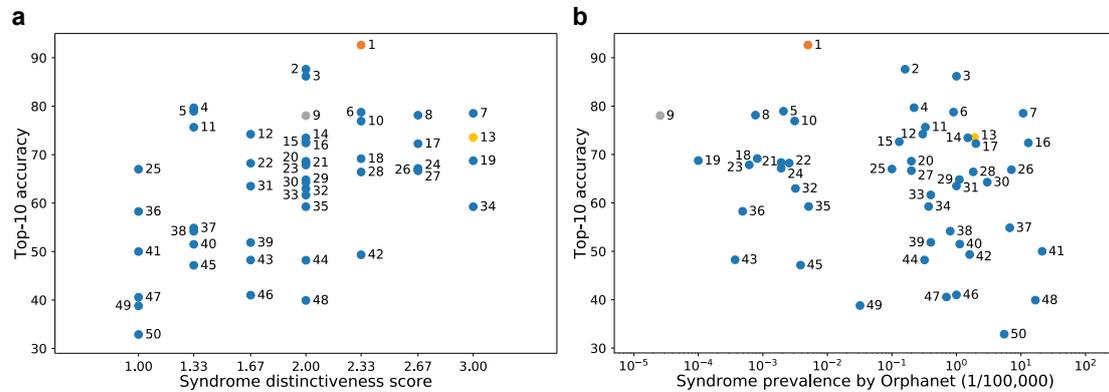


527

528 **Figure 4: Pairwise ranks of subjects with *TMEM94*.** Each label consists of family
 529 numbering and subject numbering, which are the same as in the original publication²¹.
 530 For example, F-2-7 means the seventh subject in the second family. Each column is
 531 the result of testing the image indicated at the bottom of the column. The number in
 532 the box is the rank or distance to the corresponding image in the gallery. When the
 533 rank was less than 30 and the two subjects were from different families, we added a
 534 red border to the cell. Let us take F-2-5 as an example. The sixth column starting from
 535 the left is the result of testing F-2-5, and the second row from the bottom shows that
 536 F-1-1 has a rank of 2 for F-2-5; because 2 is less than 30, a red border was added. In
 537 the third to fifth rows from the bottom are the ranks from family 2, which is the same

538 family that F-2-5 is from, so the cells do not have red borders.

539



540

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

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

543 Spearman rank correlation coefficient was 0.421 ($P = 0.002$). **b**, Distribution of top-10

544 accuracy and prevalence. The Spearman rank correlation coefficient was -0.219 ($P =$

545 0.126) The details of each syndrome can be found in Table 4 using the syndrome ID

546 shown in the figure; syndrome 9 is Schuurs-Hoeijmakers syndrome. The y-axis shows

547 the average top-10 accuracy of the experiments over those 100 iterations.

548

549 **Table 1: Performance comparison of the DeepGestalt and CASIA encoders on**
 550 **distinct, non-distinct, and target test sets.**

Test set	Model	Gallery		Test	Top-1	Top-5	Top-10	Top-30
		Images	Syndromes	images				
Distinct	Enc-DeepGestalt	20,091	296	3,083	33.56%	57.74%	68.03%	82.43%
Distinct	Enc-CASIA	20,091	296	3,083	16.65%	38.76%	51.06%	71.15%
Non-distinct	Enc-DeepGestalt	5,488.2	238.3	879.8	8.70%	22.05%	30.56%	49.84%
Non-distinct	Enc-CASIA	5,488.2	238.3	879.8	5.72%	15.87%	23.94%	42.97%
Target	Enc-DeepGestalt	2,395.3	820.4	1,186.2	11.36%	20.12%	25.25%	36.13%
Target	Enc-CASIA	2,395.3	820.4	1,186.2	7.98%	6.72%	19.00%	29.44%

551 Enc-DeepGestalt and Enc-CASIA have the same architecture. Enc-DeepGestalt training was
 552 initiated with CASIA-WebFace and further fine-tuned on photos of patients. For the top-1 to top-
 553 30 columns, the better performance in each pair is boldfaced. The numbers of images and
 554 syndromes in non-distinct and target sets are averaged over ten splits. Enc-DeepGestalt
 555 outperformed Enc-CASIA on all three types of syndromes, showing the importance of fine-
 556 tuning on patient photos for learning facial dysmorphic features.

557

558 **Table 2: Comparison of GestaltMatcher and DeepGestalt on the LMD validation**
 559 **set.**

Method	Gallery	Supported	Top-1	Top-5	Top-10	Top-30
	images	Syndromes				
DeepGestalt	-	296	54.49%	77.09%	84.52%	91.64%
GestaltMatcher	20091	296	35.91%	64.71%	74.30%	89.78%
GestaltMatcher	27826	1,362	33.74%	60.18%	70.52%	84.80%

560 The results of 323 images from LMD, validated by GestaltMatcher and DeepGestalt. We
 561 evaluated the GestaltMatcher approach on two different galleries, distinct (n = 296) and unified
 562 (n = 1,362). The best performance and the largest number of images and supported syndromes
 563 among the three conditions is boldfaced.

564

565 **Table 3: GeneMatcher validation set.**

Gene	PMID	Subject	Connected families ^a		
			Top-10	Top-30	Total
<i>BPTF</i> ³⁰	28942966	6	0	0	6
<i>CCDC47</i> ³¹	30401460	4	0	2	4
<i>CHAMP1</i> ³²	27148580	4	4	4	4
<i>CHD4</i> ³³	27616479	3	0	0	3
<i>DDX6</i> ³⁴	31422817	4	4	4	4
<i>EBF3</i> ³⁵	28017373	7	0	0	6
<i>FBXO11</i> ³⁶	30679813	17	6	9	17
<i>HNRNPK</i> ³⁷	26173930	3	3	3	3
<i>KDM3B</i> ³⁸	30929739	9	2	4	7
<i>OTUD6B</i> ³⁹	28343629	9	0	3	4
<i>PACS2</i> ⁴⁰	29656858	6	0	2	6
<i>TMEM94</i> ²¹	30526868	10	6	6	6
<i>WDR37</i> ⁴¹	31327508	4	2	2	4
<i>ZNF148</i> ⁴²	27964749	3	0	2	3
Total	-	89	27	41	77
Average	-	-	35.06%	53.25%	-

566 ^a Number of families matched by a photo from another family in the top-10 or top-30 rank.

567 For example, in the *TMEM94* study, ten out of ten images had an image from another family

568 within the top-30 rank, and all six families had at least one subject from another family in their

569 top-30 rank.

570

571 **Table 4: The 50 selected syndromes used in the random downsampling**
572 **experiment, sorted by top-10 accuracy.**

ID	Syndrome	Top 10	Score ^a	Prevalence ^b
1	Hutchinson-Gilford Progeria Syndrome; HGPS	92.62	2.33	0.005
2	Mucopolysaccharidosis Type VI; MPS6	87.63	2	0.16
3	Nijmegen Breakage Syndrome; NBS	86.18	2	1
4	Barth Syndrome; BTHS	79.67	1.33	0.22
5	Williams-Beuren Region Duplication Syndrome	78.95	1.33	0.00209
6	Crouzon Syndrome	78.76	2.33	0.9
7	Williams-Beuren Syndrome; WBS	78.53	3	10.8
8	Baraitser-Winter Syndrome	78.14	2.67	0.00077
9	Schuurs-Hoeijmakers syndrome; SHMS	78.05	2	0.00002564
10	Oculodentodigital Dysplasia	76.92	2.33	0.00312
11	Campomelic Dysplasia	75.66	1.33	0.33
12	Laron Syndrome	74.22	1.67	0.3
13	Cornelia De Lange Syndrome	73.55	3	1.9
14	Coffin-Lowry Syndrome; CLS	73.46	2	1.5
15	Pycnodysostosis	72.64	2	0.13
16	Mucopolipidosis III Alpha/beta	72.41	2	13
17	Wolf-Hirschhorn Syndrome; WHS	72.25	2.67	2
18	Renpenning Syndrome 1; RENS1	69.17	2.33	0.00082
19	Blepharophimosis, Ptosis, and Epicanthus Inversus; BPES	68.75	3	0.0001
20	Dubowitz Syndrome	68.63	2	0.2
21	Branchiooculofacial Syndrome; BOFS	68.35	2	0.00192
22	Lubs X-Linked Mental Retardation Syndrome; MRXSL	68.22	1.67	0.00256
23	Weaver Syndrome; WVS	67.86	2	0.00062
24	Hallermand-Streiff Syndrome; HSS	67.16	2.67	0.00192
25	Hyper-IgE Recurrent Infection Syndrome	67	1	0.1
26	Sotos Syndrome	66.86	2.67	7.1
27	Seckel Syndrome	66.67	2.67	0.2
28	Koolen-de Vries Syndrome; KDVS	66.42	2.33	1.82
29	Ectrodactyly, Ectodermal Dysplasia, and Cleft Lip/palate Syndrome 1; EEC1	64.81	2	1.11
30	Opitz GBBB Syndrome, Type II; GBBB2	64.29	2	3
31	Ehlers-Danlos syndrome, vascular type; EDSVASC	63.5	1.67	1
32	Simpson-Golabi-Behme Syndrome, Type 1; SGBS1	62.98	2	0.00321

33	Johanson-Blizzard Syndrome; JBS	61.64	2	0.4
34	Waardenburg Syndrome	59.26	3	0.37
35	Rothmund-Thomson Syndrome; RTS	59.26	2	0.00513
36	Mental Retardation X-Linked 102; MRX102	58.27	1	0.00049
37	Myotonic Dystrophy	54.86	1.33	6.7
38	Alagille Syndrome	54.17	1.33	0.8
39	Larsen Syndrome; LRS	51.85	1.67	0.4
40	Joubert Syndrome	51.49	1.33	1.125
41	Neurofibromatosis, Type I; NF1	50	1	21.3
42	Fetal Alcohol Syndrome; FAS	49.33	2.33	1.6
43	Filippi Syndrome; FLPIS	48.25	1.67	0.00037
44	Mucopolysaccharidosis, Type IIIA; MPS3A	48.21	2	0.32
45	Focal Dermal Hypoplasia; FDH	47.15	1.33	0.00385
46	Miller-Dieker Lissencephaly Syndrome; MDLS	41	1.67	1
47	Holt-Oram Syndrome; HOS	40.56	1	0.7
48	Trisomy 18 Syndrome	39.91	2	16.7
49	Tetrasomy 18p	38.79	1	0.03205
50	Turner Syndrome	32.88	1	5.5

573 ^a Average of the scores from three clinicians.

574 ^b Obtained from Orphanet; prevalence is per 100,000 population.

575

576 **References**

- 577 1. Ferreira, C. R. The burden of rare diseases. *Am. J. Med. Genet. Part A* **179**,
578 885–892 (2019).
- 579 2. Baird, P. A., Anderson, T. W., Newcombe, H. B. & Lowry, R. B. Genetic
580 disorders in children and young adults: A population study. *Am. J. Hum.*
581 *Genet.* **42**, 677–693 (1988).
- 582 3. Hart, T. & Hart, P. Genetic studies of craniofacial anomalies: clinical
583 implications and applications. *Orthod. Craniofac. Res.* **12**, 212–220 (2009).
- 584 4. Ferry, Q. *et al.* Diagnostically relevant facial gestalt information from ordinary
585 photos. 1–22 (2014). doi:10.7554/eLife.02020
- 586 5. Kuru, K., Niranjan, M., Tunca, Y., Osvank, E. & Azim, T. Biomedical visual
587 data analysis to build an intelligent diagnostic decision support system in
588 medical genetics. *Artif. Intell. Med.* **62**, 105–118 (2014).
- 589 6. Cerrolaza, J. J. *et al.* Identification of dysmorphic syndromes using landmark-
590 specific local texture descriptors. *2016 IEEE 13th International Symposium on*
591 *Biomedical Imaging (ISBI)* 1080–1083 (2016). doi:10.1109/ISBI.2016.7493453
- 592 7. Wang, K. & Luo, J. Detecting Visually Observable Disease Symptoms from
593 Faces. *EURASIP J. Bioinform. Syst. Biol.* **2016**, 13 (2016).
- 594 8. Dudding-Byth, T. *et al.* Computer face-matching technology using two-
595 dimensional photographs accurately matches the facial gestalt of unrelated
596 individuals with the same syndromic form of intellectual disability. *BMC*
597 *Biotechnol.* **17**, 1–9 (2017).
- 598 9. Shukla, P., Gupta, T., Saini, A., Singh, P. & Balasubramanian, R. A Deep
599 Learning Frame-Work for Recognizing Developmental Disorders. *2017 IEEE*
600 *Winter Conference on Applications of Computer Vision (WACV)* 705–714
601 (2017). doi:10.1109/WACV.2017.84
- 602 10. Liehr, T. *et al.* Next generation phenotyping in Emanuel and Pallister-Killian

- 603 syndrome using computer-aided facial dysmorphology analysis of 2D photos.
604 *Clin. Genet.* **93**, 378–381 (2018).
- 605 11. Gurovich, Y. *et al.* Identifying facial phenotypes of genetic disorders using
606 deep learning. *Nature Medicine* **25**, 60–64 (2019).
- 607 12. van der Donk, R. *et al.* Next-generation phenotyping using computer vision
608 algorithms in rare genomic neurodevelopmental disorders. *Genet. Med.* **21**,
609 1719–1725 (2019).
- 610 13. Taigman, Y., Yang, M., Ranzato, M. & Wolf, L. DeepFace: Closing the gap to
611 human-level performance in face verification. in *Proceedings of the IEEE*
612 *Computer Society Conference on Computer Vision and Pattern Recognition*
613 1701–1708 (IEEE Computer Society, 2014). doi:10.1109/CVPR.2014.220
- 614 14. Learned-Miller, E., Huang, G. B., RoyChowdhury, A., Li, H. & Hua, G. Labeled
615 faces in the wild: A survey. *Adv. Face Detect. Facial Image Anal.* 189–248
616 (2016). doi:10.1007/978-3-319-25958-1_8
- 617 15. Pantel, J. T. *et al.* Efficiency of Computer-Aided Facial Phenotyping
618 (DeepGestalt) in Individuals with and without a Genetic Syndrome: Diagnostic
619 Accuracy Study. *J. Med. Internet Res.* **22**, e19263 (2020).
- 620 16. Landrum, M. J. *et al.* ClinVar: Improving access to variant interpretations and
621 supporting evidence. *Nucleic Acids Res.* **46**, D1062–D1067 (2018).
- 622 17. McKusick, V. A. On lumpers and splitters, or the nosology of genetic disease.
623 *Perspect. Biol. Med.* **12**, 298–312 (1969).
- 624 18. Yi, D., Lei, Z., Liao, S. & Li, S. Z. Learning Face Representation from Scratch.
625 (2014).
- 626 19. Winter, R. M. & Baraitser, M. The London Dysmorphology Database. *Journal*
627 *of medical genetics* **24**, 509–510 (1987).
- 628 20. Sobreira, N., Schiettecatte, F., Valle, D. & Hamosh, A. GeneMatcher: A
629 Matching Tool for Connecting Investigators with an Interest in the Same Gene.
630 *Hum. Mutat.* **36**, 928–930 (2015).

- 631 21. Stephen, J. *et al.* Bi-allelic TMEM94 Truncating Variants Are Associated with
632 Neurodevelopmental Delay, Congenital Heart Defects, and Distinct Facial
633 Dysmorphism. *Am. J. Hum. Genet.* **103**, 948–967 (2018).
- 634 22. Alvi, M., Zisserman, A. & Nellaker, C. Turning a Blind Eye: Explicit Removal of
635 Biases and Variation from Deep Neural Network Embeddings. *Lect. Notes*
636 *Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes*
637 *Bioinformatics)* **11129 LNCS**, 556–572 (2018).
- 638 23. Lumaka, A. *et al.* Facial dysmorphism is influenced by ethnic background of
639 the patient and of the evaluator. *Clin. Genet.* **92**, 166–171 (2017).
- 640 24. Schuurs-Hoeijmakers, J. H. M. *et al.* Recurrent de novo mutations in PACS1
641 cause defective cranial-neural-crest migration and define a recognizable
642 intellectual-disability syndrome. *Am. J. Hum. Genet.* **91**, 1122–1127 (2012).
- 643 25. Van Der Maaten, L. & Hinton, G. *Visualizing Data using t-SNE.* **9**, (2008).
- 644 26. Rousseeuw, P. J. Silhouettes: a graphical aid to the interpretation and
645 validation of cluster analysis. *J. Comput. Appl. Math.* **20**, 53–65 (1987).
- 646 27. Philippakis, A. A. *et al.* The Matchmaker Exchange: A Platform for Rare
647 Disease Gene Discovery. *Hum. Mutat.* **36**, 915–921 (2015).
- 648 28. Marbach, F. *et al.* The Discovery of a LEMD2-Associated Nuclear Envelopathy
649 with Early Progeroid Appearance Suggests Advanced Applications for AI-
650 Driven Facial Phenotyping. *Am. J. Hum. Genet.* **104**, 749–757 (2019).
- 651 29. Nguengang Wakap, S. *et al.* Estimating cumulative point prevalence of rare
652 diseases: analysis of the Orphanet database. *Eur. J. Hum. Genet.* **28**, 165–
653 173 (2020).
- 654 30. Stankiewicz, P. *et al.* Haploinsufficiency of the Chromatin Remodeler BPTF
655 Causes Syndromic Developmental and Speech Delay, Postnatal
656 Microcephaly, and Dysmorphic Features. *Am. J. Hum. Genet.* **101**, 503–515
657 (2017).
- 658 31. Morimoto, M. *et al.* Bi-allelic CCDC47 Variants Cause a Disorder

- 659 Characterized by Woolly Hair, Liver Dysfunction, Dysmorphic Features, and
660 Global Developmental Delay. *Am. J. Hum. Genet.* **103**, 794–807 (2018).
- 661 32. Tanaka, A. J. *et al.* De novo pathogenic variants in CHAMP1 are associated
662 with global developmental delay, intellectual disability, and dysmorphic facial
663 features . *Mol. Case Stud.* **2**, a000661 (2016).
- 664 33. Weiss, K. *et al.* De Novo Mutations in CHD4, an ATP-Dependent Chromatin
665 Remodeler Gene, Cause an Intellectual Disability Syndrome with Distinctive
666 Dysmorphisms. *Am. J. Hum. Genet.* **99**, 934–941 (2016).
- 667 34. Balak, C. *et al.* Rare De Novo Missense Variants in RNA Helicase DDX6
668 Cause Intellectual Disability and Dysmorphic Features and Lead to P-Body
669 Defects and RNA Dysregulation. *Am. J. Hum. Genet.* **105**, 509–525 (2019).
- 670 35. Harms, F. L. *et al.* Mutations in EBF3 Disturb Transcriptional Profiles and
671 Cause Intellectual Disability, Ataxia, and Facial Dysmorphism. *Am. J. Hum.*
672 *Genet.* **100**, 117–127 (2017).
- 673 36. Jansen, S. *et al.* De novo variants in FBXO11 cause a syndromic form of
674 intellectual disability with behavioral problems and dysmorphisms. *Eur. J.*
675 *Hum. Genet.* **27**, 738–746 (2019).
- 676 37. Au, P. Y. B. *et al.* GeneMatcher Aids in the Identification of a New
677 Malformation Syndrome with Intellectual Disability, Unique Facial
678 Dysmorphisms, and Skeletal and Connective Tissue Abnormalities Caused by
679 De Novo Variants in HNRNPK. *Hum. Mutat.* **36**, 1009–1014 (2015).
- 680 38. Diets, I. J. *et al.* De Novo and Inherited Pathogenic Variants in KDM3B Cause
681 Intellectual Disability, Short Stature, and Facial Dysmorphism. *Am. J. Hum.*
682 *Genet.* **104**, 758–766 (2019).
- 683 39. Santiago-Sim, T. *et al.* Biallelic Variants in OTUD6B Cause an Intellectual
684 Disability Syndrome Associated with Seizures and Dysmorphic Features. *Am.*
685 *J. Hum. Genet.* **100**, 676–688 (2017).
- 686 40. Olson, H. E. *et al.* A Recurrent De Novo PACS2 Heterozygous Missense

- 687 Variant Causes Neonatal-Onset Developmental Epileptic Encephalopathy,
688 Facial Dysmorphism, and Cerebellar Dysgenesis. *Am. J. Hum. Genet.* **102**,
689 995–1007 (2018).
- 690 41. Kanca, O. *et al.* De Novo Variants in WDR37 Are Associated with Epilepsy,
691 Colobomas, Dysmorphism, Developmental Delay, Intellectual Disability, and
692 Cerebellar Hypoplasia. *Am. J. Hum. Genet.* **105**, 413–424 (2019).
- 693 42. Stevens, S. J. C. *et al.* Truncating de novo mutations in the Krüppel-type zinc-
694 finger gene ZNF148 in patients with corpus callosum defects, developmental
695 delay, short stature, and dysmorphisms. *Genome Med.* **8**, 131 (2016).
- 696