

1 **GestaltMatcher Database - A global reference for facial** 2 **phenotypic variability in rare human diseases**

3 Hellen Lesmann^{1,2,*}, Alexander Hustinx^{2,*}, Shahida Moosa³, Hannah Klinkhammer^{2,4},
4 Elaine Marchi⁵, Pilar Caro⁶, Ibrahim M. Abdelrazek⁷, Jean Tori Pantel^{8,9}, Merle ten
5 Hagen², Meow-Keong Thong¹⁰, Rifhan Azwani Binti Mazlan¹⁰, Sok Kun Tae¹⁰, Tom
6 Kamphans¹¹, Wolfgang Meiswinkel¹¹, Jing-Mei Li², Behnam Javanmardi², Alexej
7 Knaus², Annette Uwineza¹², Cordula Knopp¹³, Tinatin Tkemaladze^{14,15}, Miriam
8 Elbracht¹³, Larissa Mattern¹³, Rami Abou Jamra¹⁶, Clara Velmans¹⁷, Vincent
9 Strehlow¹⁶, Maureen Jacob¹⁸, Angela Peron^{19,20}, Cristina Dias^{21,22,23,24}, Beatriz
10 Carvalho Nunes²⁵, Thainá Vilella²⁵, Isabel Furquim Pinheiro²⁶, Chong Ae Kim²⁶,
11 Maria Isabel Melaragno²⁵, Hannah Weiland², Sophia Kaptain², Karolina
12 Chwiałkowska^{27,28}, Mirosław Kwasniewski^{28,27}, Ramy Saad^{22,29}, Sarah Wiethoff³⁰,
13 Himanshu Goel³¹, Clara Tang³², Anna Hau³³, Tahsin Stefan Barakat³⁴, Przemysław
14 Panek³⁵, Amira Nabil⁷, Julia Suh¹³, Frederik Braun³⁶, Israel Gomy³⁷, Luisa
15 Averdunk³⁸, Ekanem Ekure³⁹, Gaber Bergant⁴⁰, Borut Peterlin⁴¹, Claudio Graziano⁴²,
16 Nagwa Gaboon^{43,44}, Moisés Fiesco-Roa^{45,46}, Alessandro Mauro Spinelli⁴⁷, Nina-
17 Maria Wilpert^{48,49,50}, Prasit Phowthongkum^{51,52}, Nergis Güzel¹³, Tobias B. Haack⁵³,
18 Rana Bitar^{54,55}, Andreas Tzschach⁵⁶, Agusti Rodriguez-Palmero⁵⁷, Theresa Brunet¹⁸,
19 Sabine Rudnik-Schöneborn⁵⁸, Silvina Noemi Contreras-Capetillo⁵⁹, Ava Oberlack¹⁸,
20 Carole Samango-Sprouse^{60,61,62}, Teresa Sadeghin⁶³, Margaret Olaya⁶³, Konrad
21 Platzer¹⁶, Artem Borovikov⁶⁴, Franziska Schnabel¹⁶, Lara Heuft¹⁶, Vera Herrmann¹⁶,
22 Renske Oegema⁶⁵, Nour Elkhateeb⁶⁶, Sheetal Kumar¹, Katalin Komlosi⁵⁶,
23 Khoushoua Mohamed⁷, Silvia Kalantari⁶⁷, Fabio Sirchia^{67,68}, Antonio F. Martinez-
24 Monseny⁶⁹, Matthias Höller⁵⁶, Louiza Toutouna⁵⁶, Amal Mohamed⁷, Amaia Las-
25 Aranzasti^{70,71}, John A. Sayer^{72,73}, Nadja Ehmke⁷⁴, Magdalena Danyel⁷⁴, Henrike
26 Sczakiel⁷⁴, Sarina Schwartzmann⁷⁴, Felix Boschann⁷⁴, Max Zhao⁷⁴, Ronja Adam⁷⁴,
27 Lara Einicke⁷⁴, Denise Horn⁷⁴, Kee Seang Chew⁷⁵, Choy Chen KAM⁷⁵, Miray
28 Karakoyun⁷⁶, Ben Pode-Shakked^{77,78}, Aviva Eliyahu^{79,80,81}, Rachel Rock^{82,83}, Teresa
29 Carrion⁸⁴, Odelia Chorin⁸⁵, Yuri A. Zarate^{86,87}, Marcelo Martinez Conti⁸⁸, Mert
30 Karakaya¹⁷, Moon Ley Tung^{89,90}, Bharatendu Chandra^{89,90}, Arjan Bouman³⁴, Aime
31 Lumaka⁹¹, Naveed Wasif^{92,93}, Marwan Shinawi⁹⁴, Patrick R. Blackburn⁹⁵, Tianyun
32 Wang^{96,97,98}, Tim Niehues⁹⁹, Axel Schmidt¹, Regina Rita Roth¹⁰⁰, Dagmar
33 Wieczorek¹⁰⁰, Ping Hu¹⁰¹, Rebekah L. Waikel¹⁰¹, Suzanna E. Ledgister Hanchard¹⁰¹,
34 Gehad Elmakkawy⁷, Sylvia Safwat⁷, Frédéric Ebstein^{102,103}, Elke Krüger¹⁰⁴,
35 Sébastien Küry^{102,103}, Stéphane Bézieau^{102,103}, Annabelle Arlt², Eric Olinger¹⁰⁵, Felix
36 Marbach⁶, Dong Li¹⁰⁶, Lucie Dupuis¹⁰⁷, Roberto Mendoza-Londono¹⁰⁷, Sofia
37 Douzgou Houge¹⁰⁸, Denisa Weis¹⁰⁹, Brian Hon-Yin Chung^{110,111}, Christopher C.Y.
38 Mak¹¹¹, Hülya Kayserili¹¹², Nursel Elcioglu¹¹³, Ayca Aykut¹¹⁴, Peli Özlem Şimşek-
39 Kiper¹¹⁵, Nina Bögershausen¹¹⁶, Bernd Wollnik^{116,117,118}, Heidi Beate Bentzen^{119,120},
40 Ingo Kurth¹³, Christian Netzer¹⁷, Aleksandra Jezela-Stanek³⁵, Koen Devriendt¹²¹,
41 Karen W. Gripp¹²², Martin Mücke^{8,9}, Alain Verloes¹²³, Christian P. Schaaf⁶,

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

42 Christoffer Nellåker¹²⁴, Benjamin D. Solomon¹⁰¹, Markus M. Nöthen¹, Ebtesam
43 Abdalla⁷, Gholson J. Lyon^{125,126,127}, Peter M. Krawitz², Tzung-Chien Hsieh^{2,#}

44
45

46 ¹Institute of Human Genetics, University of Bonn, Bonn, NRW, Germany, ²Institute
47 for Genomic Statistics and Bioinformatics, University of Bonn, Bonn, NRW,
48 Germany, ³Division of Molecular Biology and Human Genetics, Stellenbosch
49 University and Medical Genetics, Tygerberg Hospital, Stellenbosch, South Africa,
50 ⁴Institute for Medical Biometry, Informatics and Epidemiology, University of Bonn,
51 Bonn, NRW, Germany, ⁵New York State Institute for Basic Research in
52 Developmental Disabilities, New York State, Albany, New York, USA, ⁶Institute of
53 Human Genetics, Heidelberg University, Heidelberg, Baden-Württemberg, Germany,
54 ⁷Department of Human Genetics, Medical Research Institute, Alexandria University,
55 Alexandria, Alexandria, Egypt, ⁸Institute for Digitalization and General Medicine,
56 University Hospital RWTH Aachen, Aachen, NRW, Germany, ⁹Centre for Rare
57 Diseases Aachen (ZSEA), University Hospital RWTH Aachen, Aachen, NRW,
58 Germany, ¹⁰Department of Paediatrics, Faculty of Medicine, University of Malaya,
59 50603 Kuala Lumpur, Malaysia, ¹¹GeneTalk GmbH, Bonn, NRW, Germany,
60 ¹²College of Medicine and Health Sciences, University of Rwanda, and University
61 Teaching Hospital of Kigali, Kigali, Rwanda, ¹³Institute for Human Genetics and
62 Genomic Medicine, Medical Faculty, RWTH Aachen University, Aachen, NRW,
63 Germany, ¹⁴Department of Molecular and Medical Genetics, Tbilisi State Medical
64 University, Tbilisi, Georgia, ¹⁵Givi Zhvania Pediatric Academic Clinic, Tbilisi State
65 Medical University, Georgia, ¹⁶Institute of Human Genetics, University of Leipzig
66 Medical Center, Leipzig, Germany, ¹⁷Institute of Human Genetics, University of
67 Cologne, Faculty of Medicine and University Hospital Cologne, Cologne, NRW,
68 Germany, ¹⁸Institute of Human Genetics, Klinikum rechts der Isar, Technical
69 University of Munich, School of Medicine and Health, Munich, Germany, ¹⁹Medical
70 Genetics, Meyer Children's Hospital IRCCS, Firenze, Italy, ²⁰Department of
71 Experimental and Clinical Biomedical Sciences "Mario Serio", Università degli Studi
72 di Firenze, Italy, ²¹Department of Medical Genetics, Guy's and St. Thomas' NHS
73 Foundation Trust, London, UK, ²²North East Thames Regional Genetics Service,
74 Great Ormond Street Hospital for Children, Great Ormond Street, London, UK,
75 ²³Neural Stem Cell Biology Laboratory, The Francis Crick Institute, UK, ²⁴Department
76 of Medical & Molecular Genetics, School of Basic and Medical Biosciences, Faculty
77 of Life Sciences & Medicine, King's College London, UK, ²⁵Genetics Division,
78 Department of Morphology and Genetics, Universidade Federal de São Paulo, São
79 Paulo, Brazil, ²⁶Genetics Unit, Instituto da Criança, Universidade de São Paulo, São
80 Paulo, Brazil, ²⁷Centre for Bioinformatics and Data Analysis, Medical University of
81 Bialystok, Bialystok, Poland, ²⁸IMAGENE.ME SA, Bialystok, Poland, ²⁹Department of
82 Genetics and Genomic Medicine, UCL Institute of Child Health, London UK,
83 ³⁰Department of Neurology with Institute of Translational Neurology, University
84 Hospital Münster, Münster, NRW, Germany, ³¹School of Medicine and Public Health,
85 University of Newcastle, Callaghan NSW, Australia, ³²Kabuki Syndrome Foundation,

86 Northbrook, IL, USA, ³³Hunter Genetics, Hunter New England Health Service,
87 Newcastle, Australia, ³⁴Department of Clinical Genetics, Erasmus MC University
88 Medical Center, Rotterdam, The Netherlands, ³⁵Department of Genetics and Clinical
89 Immunology, National Institute of Tuberculosis and Lung Diseases, Warsaw, Poland,
90 ³⁶Institute of Human Genetics, University Hospital Essen, University Duisburg-Essen,
91 Essen, NRW, Germany, ³⁷Department of Genetics, Faculdade de Medicina de
92 Ribeirão Preto, Universidade de São Paulo, Sao Paulo, Brazil, ³⁸Department of
93 General Pediatrics and Neonatology, University Children's Hospital, Heinrich-Heine-
94 University, Medical Faculty, Düsseldorf, Germany, ³⁹Department of Paediatrics,
95 College of Medicine, University of Lagos, Lagos, Nigeria, ⁴⁰Clinical Institute of
96 Genomic Medicine, University Medical Centre Ljubljana, Ljubljana, Slovenia,
97 ⁴¹Clinical Institute of Genomic Medicine, University Medical Centre Ljubljana,
98 ⁴²Medical Genetics Unit, Ausl Romagna, Cesena, Italy, ⁴³Medical Genetics Center,
99 Faculty of Medicine, Ain Shams University, Cairo, Egypt, ⁴⁴Medical Genetics
100 Department, Armed Forces College of Medicine, Cairo, Egypt, ⁴⁵Programa de
101 Maestría y Doctorado en Ciencias Médicas, Odontológicas y de la Salud,
102 Universidad Nacional Autónoma de México, México City, Mexico, ⁴⁶Laboratorio de
103 Citogenética, Instituto Nacional de Pediatría, México City, Mexico, ⁴⁷Institute for
104 Maternal and Child Health, IRCCS Burlo Garofolo, Trieste, Italy, ⁴⁸NeuroCure Cluster
105 of Excellence; Charité–Universitätsmedizin Berlin, corporate member of Freie
106 Universität Berlin and Humboldt-Universität zu Berlin, D-10117 Berlin, Germany,
107 ⁴⁹Department of Neuropediatrics, Charité–Universitätsmedizin Berlin, corporate
108 member of Freie Universität Berlin and Humboldt-Universität zu Berlin, D-13353
109 Berlin, Germany, ⁵⁰Berlin Institute of Health at Charité - Universitätsmedizin Berlin,
110 BIH Biomedical Innovation Academy, BIH Charité Junior Clinician Scientist Program,
111 D-10117 Berlin, German, ⁵¹Excellence Center for Genomics and Precision Medicine,
112 King Chulalongkorn Memorial Hospital, the Thai Red Cross Society, Bangkok,
113 Thailand, ⁵²Division of Medical Genetics and Genomics, Department of Medicine,
114 Faculty of Medicine, Chulalongkorn University, Bangkok, Thailand, ⁵³Institute of
115 Medical Genetics and Applied Genomics, University of Tübingen, Tübingen,
116 Germany, ⁵⁴Pediatric Gastroenterology Department, Sheikh Khalifa Medical City,
117 Abu Dhabi, United Arab Emirates, ⁵⁵Khalifa University, Abu Dhabi, United Arab
118 Emirates, ⁵⁶Institute of Human Genetics, Medical Center - University of Freiburg,
119 Faculty of Medicine, University of Freiburg, Freiburg, Germany, ⁵⁷Paediatric
120 Neurology Unit, Department of Pediatrics, Hospital Universitari Germans Trias i
121 Pujol, Universitat Autònoma de Barcelona, Barcelona, Spain, ⁵⁸Institute of Human
122 Genetics, Medical University Innsbruck, Innsbruck, Austria, ⁵⁹Universidad Autónoma
123 de Yucatan (University Autonomus of Yucatan), Merida, Yucatan, Mexico,
124 ⁶⁰Department of Pediatrics, George Washington University, 2121 I St. NW,
125 Washington D.C. 2005, ⁶¹Department of Human and Molecular Genetics, Florida
126 International University, 11200 SW 8th Street, AHC2 Miami, Florida 22199,
127 ⁶²Department of Research, The Focus Foundation, 820 W. Central Ave. #190,
128 Davidsonville, MD 21035, ⁶³Department of Research, The Focus Foundation, 2772
129 Rutland Road P.O. Box 190, Davidsonville, MD 21035, ⁶⁴Research Centre for

130 Medical Genetics (RCMG), Moscow, Russia, ⁶⁵Department of Genetics, University
131 Medical Center Utrecht, Utrecht University, Utrecht, The Netherlands, ⁶⁶Department
132 of Clinical Genetics, Cambridge University Hospitals NHS Foundation Trust,
133 Cambridge, UK, ⁶⁷Department of Molecular Medicine, University of Pavia, Pavia,
134 Italy, ⁶⁸Medical Genetics Unit, IRCCS San Matteo Foundation, Pavia, Italy,
135 ⁶⁹Department of Clinical Genetics, SJD Barcelona Children's Hospital, Esplugues del
136 Llobregat (Barcelona), Spain, ⁷⁰Medicine Genetics Group, Vall d'Hebron Institut de
137 Recerca (VHIR), Vall d'Hebron Barcelona Hospital Campus, Vall d'Hebron Hospital
138 Universitari, Barcelona, Catalunya, Spain, ⁷¹Department of Clinical and Molecular
139 Genetics, Vall d'Hebron Barcelona Hospital Campus, Vall d'Hebron Hospital
140 Universitari, Barcelona, Catalunya, Spain, ⁷²Biosciences Institute, Newcastle
141 University, Central Parkway, Newcastle upon Tyne, UK, ⁷³Renal Services, The
142 Newcastle Upon Tyne NHS Hospitals Foundation Trust, Freeman Road, Newcastle
143 Upon Tyne, UK, ⁷⁴Institute of Medical Genetics and Human Genetics, Charité-
144 Universitätsmedizin Berlin, Humboldt-Universität zu Berlin and Berlin Institute of
145 Health, Berlin, Germany, ⁷⁵Department of Paediatrics, Faculty of Medicine,
146 University Malaya, 59100 Kuala Lumpur, Malaysia, ⁷⁶Ege University, Faculty of
147 Medicine, Department of Pediatric Gastroenterology Hepatology and Nutrition, Izmir,
148 Turkey, ⁷⁷The Institute of Rare Diseases, Edmond and Lily Safra Children's Hospital,
149 Sheba Medical Center, Ramat Gan, Israel, ⁷⁸The faculty of Medical and Health
150 Sciences, Tel-Aviv University, Tel-Aviv, Israel, ⁷⁹The Danek Gertner Institute of
151 Human Genetics, Sheba Medical Center, Tel-Hashomer, Israel, ⁸⁰Sackler Faculty of
152 Medicine, Tel-Aviv University, Tel-Aviv, Israel, ⁸¹The Mina and Everard Goodman
153 Faculty of Life Sciences, Bar-Ilan University, Ramat Gan, Israel, ⁸²Metabolic
154 Diseases Clinic, Edmond and Lily Safra Children's Hospital, Sheba Medical Center,
155 ⁸³National Newborn Screening Program, Public Health Services, Ministry of Health
156 Tel-Hashomer, Israel, ⁸⁴Rare diseases Unit, Pediatric Department, Hospital
157 Universitari Son Espases, Palma de Mallorca, Spain, ⁸⁵The Institute of Rare
158 Diseases, Edmond and Lily Safra Children's Hospital, Sheba Medical Center, Tel-
159 Hashomer, Israel, ⁸⁶Department of Pediatrics, Section of Genetics and Metabolism,
160 University of Arkansas for Medical Sciences and Arkansas Children's Hospital, Little
161 Rock, AR, USA, ⁸⁷Division of Genetics and Metabolism, University of Kentucky,
162 Lexington, KY, USA, ⁸⁸Project Director AI for Health at Foundation 29, Foundation
163 29, Madrid, Spain, ⁸⁹University of Iowa Roy J and Lucille A Carver College of
164 Medicine, Iowa City, IA 52242, USA, ⁹⁰Division of Medical Genetics and Genomics,
165 Stead Family Department of Pediatrics, University of Iowa Hospitals and Clinics,
166 Iowa City, IA 52242, USA, ⁹¹Center for Human Genetics, Faculty of Medicine,
167 University of Kinshasa, Kinshasa, DR Congo, ⁹²Institute of Human Genetics,
168 University of Ulm, Ulm, Baden-Württemberg, Germany, ⁹³University Hospital
169 Schleswig-Holstein, Campus Kiel, Kiel, Germany, ⁹⁴Division of Genetics and
170 Genomic Medicine, Department of Pediatrics, Washington University School of
171 Medicine, St. Louis, MO, USA, ⁹⁵Department of Pathology, St. Jude Children's
172 Research Hospital, Memphis, Tennessee 38105, USA, ⁹⁶Department of Medical
173 Genetics, Center for Medical Genetics, Peking University Health Science Center,

174 Beijing 100191, China, ⁹⁷Neuroscience Research Institute, Peking University; Key
175 Laboratory for Neuroscience, Ministry of Education of China & National Health
176 Commission of China, Beijing 100191, China, ⁹⁸Autism Research Center, Peking
177 University Health Science Center, Beijing 100191, China, ⁹⁹Department of Pediatrics,
178 Helios Klinik Krefeld, Krefeld 47805, Germany, ¹⁰⁰Institute of Human Genetics,
179 Medical Faculty, University Hospital Düsseldorf, Heinrich Heine University
180 Düsseldorf, Germany, ¹⁰¹Medical Genomics Unit, Medical Genetics Branch, National
181 Human Genome Research Institute, Bethesda, USA, ¹⁰²Nantes Université, CHU
182 Nantes, CNRS, INSERM, l'institut du thorax, F-44000 Nantes, France, ¹⁰³Nantes
183 Université, CHU Nantes, Service de Génétique Médicale, F-44000 Nantes, France,
184 ¹⁰⁴Institute for Medical Biochemistry and Molecular Biology, University of Greifswald,
185 Greifswald, Greifswald, Germany, ¹⁰⁵Center for Human Genetics, Cliniques
186 Universitaires Saint-Luc, Brussels, Belgium, ¹⁰⁶Division of Human Genetics,
187 Children's Hospital of Philadelphia, Philadelphia, Pennsylvania, USA, ¹⁰⁷Department
188 to Paediatrics, Division of Clinical and Metabolic Genetics, The Hospital of Sick
189 Children, Toronto, Ontario, Canada, ¹⁰⁸Department of Medical Genetics, Haukeland
190 University Hospital, Bergen, Norway, ¹⁰⁹Institute for Medical Genetics, Kepler
191 University Hospital, Linz, Austria, ¹¹⁰Hong Kong Genome Institute, Hong Kong,
192 China, ¹¹¹Department of Paediatrics and Adolescent Medicine, The University of
193 Hong Kong, Hong Kong, China, ¹¹²Medical Genetics Department, Koç University
194 School of Medicine (KUSoM), 34010, Istanbul, Türkiye, ¹¹³Department of Pediatric
195 Genetics, Marmara University School of Medicine, Istanbul, Türkiye, ¹¹⁴Department
196 of Medical Genetics, Ege University Faculty of Medicine, Izmir, Türkiye, ¹¹⁵Hacettepe
197 University Faculty of Medicine, Department of Pediatric Genetics, Ankara, Türkiye,
198 ¹¹⁶Institut of Human Genetics, University Medical Center Göttingen, Göttingen,
199 Germany, ¹¹⁷Cluster of Excellence "Multiscale Bioimaging: from Molecular Machines
200 to Networks of Excitable Cells" (MBExC), University of Göttingen, Göttingen,
201 Germany, ¹¹⁸German Center for Cardiovascular Research (DZHK), Partner Site
202 Göttingen, Göttingen, Germany, ¹¹⁹Centre for Medical Ethics, Faculty of Medicine,
203 University of Oslo, Oslo, Norway, ¹²⁰Cancer Registry of Norway, Norwegian Institute
204 of Public Health, Oslo, Norway, ¹²¹Center for Human Genetics, KU Leuven, Leuven,
205 Belgium, ¹²²Division of Medical Genetics, A.I. du Pont Hospital for Children/Nemours,
206 USA, Wilmington, Delaware, USA, ¹²³Department of Clinical Genetics, Robert-Debré
207 Hospital, Paris, France, ¹²⁴Big Data Institute, Li Ka Shing Centre for Health
208 Information and Discovery, Nuffield Department of Women's & Reproductive Health,
209 University of Oxford, Oxford, UK, ¹²⁵Department of Human Genetics, New York State
210 Institute for Basic Research in Developmental Disabilities, Staten Island, New York,
211 United States of America, ¹²⁶George A. Jervis Clinic, New York State Institute for
212 Basic Research in Developmental Disabilities, Staten Island, New York, United
213 States of America, ¹²⁷Biology PhD Program, The Graduate Center, The City
214 University of New York, New York, United States of America
215

216 *These authors contributed equally

217 #Corresponding author

218 Abstract

219 The most important factor that complicates the work of dysmorphologists is the
220 significant phenotypic variability of the human face. Next-Generation Phenotyping
221 (NGP) tools that assist clinicians with recognizing characteristic syndromic patterns
222 are particularly challenged when confronted with patients from populations different
223 from their training data. To that end, we systematically analyzed the impact of genetic
224 ancestry on facial dysmorphism. For that purpose, we established the GestaltMatcher
225 Database (GMDB) as a reference dataset for medical images of patients with rare
226 genetic disorders from around the world. We collected 10,980 frontal facial images –
227 more than a quarter previously unpublished - from 8,346 patients, representing 581
228 rare disorders. Although the predominant ancestry is still European (67%), data from
229 underrepresented populations have been increased considerably via global
230 collaborations (19% Asian and 7% African). This includes previously unpublished
231 reports for more than 40% of the African patients. The NGP analysis on this diverse
232 dataset revealed characteristic performance differences depending on the
233 composition of training and test sets corresponding to genetic relatedness. For clinical
234 use of NGP, incorporating non-European patients resulted in a profound enhancement
235 of GestaltMatcher performance. The top-5 accuracy rate increased by +11.29%.
236 Importantly, this improvement in delineating the correct disorder from a facial portrait
237 was achieved without decreasing the performance on European patients. By design,
238 GMDB complies with the FAIR principles by rendering the curated medical data
239 findable, accessible, interoperable, and reusable. This means GMDB can also serve
240 as data for training and benchmarking. In summary, our study on facial dysmorphism
241 on a global sample revealed a considerable cross ancestral phenotypic variability
242 confounding NGP that should be counteracted by international efforts for increasing
243 data diversity. GMDB will serve as a vital reference database for clinicians and a
244 transparent training set for advancing NGP technology.

245 Introduction

246 Facial dysmorphism is one of the most complex and informative clinical features in
247 syndromic disorders, and is therefore often crucial in terms of establishing a diagnosis

248 in rare genetic diseases^{1,2}. However, the recognition of dysmorphic patterns, is a
249 challenging endeavour, and relies on the skills, knowledge, and experience of the
250 examiner. In certain syndromes, in particular those that are ultra-rare, variability in
251 facial features can pose challenges even for highly experienced clinicians³. Facial
252 features can also vary according to sex, age, and ancestry, which further complicates
253 the recognition of a specific dysmorphic pattern^{4–6}.

254 Ancestry plays a particularly significant role since considerable inter-ancestral
255 variability exists in facial gestalt⁷. Thus, facial features that are common in certain
256 ancestral groups may be considered dysmorphic in others. For example, while
257 upslanting palpebral fissures are common in healthy Asians, they may be perceived
258 as dysmorphic in other populations⁸. Previous studies have also highlighted
259 differences in facial gestalt between different ancestries in common dysmorphic
260 genetic syndromes such as Down Syndrome, 22q11.2 deletion syndrome, Noonan
261 syndrome, and Williams–Beuren syndrome^{4,9,10}. Furthermore, Lumaka et al. have
262 demonstrated that this variability can influence the assessor, with European clinicians
263 failing to recognize dysmorphic features in individuals of African ancestry¹¹. This is a
264 growing problem as globalization and migration increasingly blur ancestral and cultural
265 boundaries, and geography is no longer a key determining factor in mating patterns¹².
266 Hence, in diverse populations, such as those with admixed ancestries, the challenge
267 of accurately diagnosing rare diseases becomes even more pronounced since new
268 phenotypes can evolve via admixture¹³.

269 Ancestry also has a significant impact on the detection of rare dysmorphic disorders
270 via artificial intelligence (AI)¹¹ because in most healthcare datasets, non-European
271 ancestries are underrepresented¹⁴. Many next-generation phenotyping (NGP)
272 approaches that predict disorders on the basis of facial image analysis, such as
273 GestaltMatcher¹⁵, have demonstrated high accuracy in patients from the ancestries in
274 which they were predominantly trained and validated, i.e., European and North
275 American^{15–19}.

276 Since the significantly higher birth rates in non-European regions account for 80% of
277 the global population and 90% of all annual births (Figure 1a)²⁰, action is required to
278 include non-European patients currently considered to be underrepresented. So far,
279 few studies exist about the performance of NGP tools where the ancestry composition

280 of individuals in the training and test set differs. Literature suggests that AIs trained on
281 individuals of European ancestry perform better on a test set of Asian rather than
282 African ancestry²¹⁻²⁴ that may be explained by their closer genetic relatedness²⁵. This
283 raises the question of whether AIs need to be trained for different ancestries or whether
284 a similar performance can be achieved by sufficiently increasing the ancestral diversity
285 in the joint training set. The latter is indicated by a study conducted on Down syndrome
286 patients of African ancestry¹¹. However, comparing these studies is difficult since they
287 were not performed on data compliant with FAIR principles that are findable,
288 accessible, interoperable, and reusable, meaning the results cannot be

289 The motivation of our work is therefore threefold: 1) scientific, because we wanted to
290 study the effect of inter- and intra-ancestral phenotypic variability on NGP, such as
291 GestaltMatcher, in a systematic manner; 2) clinical, because more diverse training
292 data can presumably increase the performance of NGP on non-European ancestries;
293 and 3) societal, because so far underrepresented populations would benefit from
294 potential performance improvements.

295 To achieve these goals, we aimed for a FAIR database with an increased number of
296 patients of non-European ancestry with respect to comparable databases^{20,26,27}.
297 Therefore, we established the GestaltMatcher Database (GMDB) as a community-
298 driven online framework that facilitates acquiring patient consent and incentivizes data
299 sharing, acknowledging contributions from clinician-scientists as citeable micro-
300 publications (Figure 2)²⁸⁻³¹. Through this framework, we established global
301 collaborations, enabling the collection of a wide range of data from various ancestries.

302 GMDB is the first database for medical imaging data of patients with rare genetic
303 disorders from diverse ancestries that is compliant with the FAIR principles³². By its
304 machine-readable design, GMDB also enables systematic analyses of the influence
305 of genetic background on NGP performance, which we will report in this study.

306 Results:

307 Overview of FAIR data in GMDB

308 Retrospective data from curated publications, along with data provided by clinicians or
309 patients, were made available as FAIR cases in the GMDB (Figure 3, Supplementary
310 Figures 1 and 2)³³. At the time of the data freeze for this paper on April 6th 2024, we

311 curated the GMDB-FAIR dataset consisting of 10,980 portrait images (Supplementary
312 Figure 3) of 8,346 patients with 581 genetic disorders, including patients curated from
313 2,224 scientific publications. 2,312 unpublished images were contributed by 138
314 clinicians from 106 institutions (indicated by location markers in Figure 1a), including
315 novel cases from GMDB micro-publications (micro-publication section in
316 Supplementary Note). For the portrait data, which is the scope of this study, in terms
317 of sex, the data distribution is relatively balanced (Figure 4a). However, age is biased
318 toward patients aged below 10 years (Figure 4b). Figure 4c shows a two-dimensional
319 representation of Human Phenotype Ontology³⁴ (HPO)-defined symptom groups in
320 GMDB via Uniform Manifold Approximation and Projection (UMAP). While GMDB
321 incorporates cases from all HPO-defined symptom groups across the disease
322 landscape, the HPO-defined symptom group ‘facial dysmorphism’ is enriched in
323 GMDB. Since each individual can be attributed to several HPO-defined symptom
324 groups according to their features, facial dysmorphism was also present in the other
325 HPO-defined symptom groups, as shown in the heatmap.

326 **Underrepresented populations benefited from micro-publication case reports in** 327 **GMDB**

328 Through our international collaborations (Figure 1a), the representation of non-
329 European ancestral groups is 19% for Asian, 7% for African, and 7% for Others. 67%
330 comprises individuals of European descent (Figure 1b). Moreover, the ancestry
331 distribution varies among different disorders. Some disorders, such as Williams-
332 Beuren syndrome, Hyperphosphatasia with impaired intellectual development
333 syndrome, and Cohen syndrome, have relatively diverse and balanced ancestral
334 distributions (Supplementary Figure 4).

335 Notably, the proportion of African ancestry was strongly increased by means of GMDB
336 micro-publications which account for 40% of the individuals with African ancestry
337 (Figure 4d). In terms of specific sub-ancestries (Figure 4e), more than 80% of cases
338 with sub-Saharan ancestry and over 20% of cases with North African, Native American,
339 and Latin American ancestries were obtained through GMDB micro-publications.

340 **Performance disparities in underrepresented populations**

341 We analyzed the performance of GestaltMatcher on the test set of 882 images of 275
342 disorders with different ancestries that have not been used for the training of
343 GestaltMatcher. Performance is measured as a top-k accuracy (as described in
344 Methods). We report the top-1 to top-30 accuracies in Table 1. When considering top-
345 1 accuracy, the 'Others' group demonstrated the highest performance at 73.91%,
346 followed by the African group at 62.07%, the Asian group at 53.54%, and the European
347 group at 55.45%. The African group achieved the highest top-5 accuracy (82.76%),
348 the Asian group attained the highest top-10 accuracy (85.04%), while the European
349 group only achieved 75.14% and 82.60% for top-5 and top-10 accuracies, respectively.
350 However, the European group contains more than 50% of the testing images (523 out
351 of 882), covering many more disorders than the other ancestry groups. That includes
352 ultra-rare disorders known to achieve lower performances¹⁹.

353 To fairly compare the European group to another non-European ancestry, we only
354 looked at the disorders that were present in both ancestry groups. In

355

356

357 **Table 2**, when comparing the African and European groups on the six overlapping
358 disorders, the European group outperformed the African group by achieving +16.96%
359 top-1 accuracy and +11.17% top-10 accuracy. The European group also exhibited
360 higher accuracies compared to the Asian group, with a top-1 accuracy of +6.92% and
361 a top-10 accuracy of +4.15%. However, the European and 'Others' groups achieved
362 relatively comparable results. The 'Others' group had a higher top-1 accuracy, while
363 the European group performed better on the top-10 accuracy.

364 We further reported the performance of sex and age groups in Table 1. The distribution
365 of testing images was relatively balanced across different groups, and no significant
366 performance gap was observed between males and females. However, the under-
367 one-year-old group exhibited the lowest performance, while the five- to ten-year-old
368 group demonstrated notably higher top-5 and top-10 accuracies.

369 **Diverse ancestry data enhance prediction accuracy for underrepresented** 370 **populations**

371 To investigate the impact of incorporating ancestry-diverse data on the overall
372 performance of GestaltMatcher across ancestries, we designed two sets of ancestry
373 analysis experiments. First, we investigated the expansion of the training set of
374 GestaltMatcher (as described in Methods), including either European only (EU + EU*)
375 or European and non-European (EU + non-EU) patients. We measured a top-1
376 accuracy averaged over all ancestral groups of 49.65% for the European only training
377 set (EU + EU*) and 66.90% for the diverse training set (EU + non-EU) (Figure 5a).
378 Similarly, top-5 accuracy of the European training set was 69.95%, and when we
379 trained on the diverse set, the top-5 accuracy increased to 81.24%. Notably, the
380 evaluation performance on images of patients with European ancestry showed only a
381 marginal performance dropdown. Specifically, the top-1 accuracy decreased by 3.82%
382 and the top-5 accuracy by 3.61% when the dataset was augmented with 50% more
383 non-European images. Meanwhile, the top-1 and top-5 performance increased notably
384 for almost every other ancestral group. Figure 5a and Table 3 show further per-
385 ancestry performances.

386 The training of GestaltMatcher results in a clinical face phenotype space that can be
387 populated by additional cases, which we refer to as the gallery set (as described in
388 Methods). We next investigated the influence of expanding the gallery with ancestry-
389 diverse data by gradually raising the proportion of included non-European data from
390 10% to 100%. Figure 5b shows that the top-1 accuracy of the non-European groups
391 was clearly increased when we added more non-European data in the gallery.
392 However, the top-1 accuracy of the European group did not change even when we
393 added 100% of the non-European data into the gallery.

394 **GMDB-FAIR dataset drives the advancement of NGP technology**

395 GMDB-FAIR dataset is the first dataset that can be shared with the research
396 community to train and benchmark their NGP approaches. After the first publication of
397 the GestaltMatcher approach in 2022, for which we initially started the collection of our
398 FAIR data, many researchers have utilized GMDB-FAIR to develop different NGP
399 approaches. Hustinx et al.¹⁹, Sumer et al.³⁵, and Campbell et al.³⁶ improved the
400 prediction accuracy of their models significantly by utilizing different loss functions,

401 network architectures, and data augmentation. Recently, Wu et al. proposed
402 combining a large language model with facial image analysis to streamline the rare
403 disorder diagnosis³⁷. Furthermore, running facial analysis with an on-premise solution
404 is possible using the FAIR data set to further prioritize genomic variants³⁸.

405 Moreover, the GMDB-FAIR dataset can be taken as a validatable control cohort to
406 facilitate the delineation of the facial phenotype of disorders. GestaltMatcher can
407 detect clusters and assess whether, for example, cases with an identical variant or
408 pathogenic variants in the same gene share a similar facial phenotype. For example,
409 Ebstein et al. showed that facial dysmorphism was heterogeneous among the entire
410 *PSMC3* patient cohort, but facial similarities were found in patients sharing the same
411 pathogenic variants³⁹. To date, 15 publications have analyzed the facial phenotype of
412 the cohort with the GMDB-FAIR dataset and GestaltMatcher³⁹⁻⁵³. All results can be
413 reproduced in the research platform of GMDB, which we introduce in the Methods
414 section (Figure 2c, Figure 3c and Supplementary Note).

415 Discussion

416 GMDB is a modern, searchable reference and publication medium encompassing
417 diverse populations that is designed for both clinicians and computer scientists
418 engaged in NGP development. The ultimate goal of this study is to drive research in
419 rare genetic disorders to understand the phenotypic variability among ancestries
420 systematically and improve support for underrepresented populations.

421 GMDB stands out as the sole database compliant with FAIR principles, distinguished
422 by its extensive collection of facial images covering diverse populations. This was
423 mainly possible through the contributions and crowd-sourced annotations by our
424 global collaborators. To increase motivation for data submission in the future, every
425 case in the database has the potential to become a citable micro-publication with a
426 Digital Object Identifier (DOI)⁵⁴. Furthermore, future micro-publications could be
427 indexed in reputable scientific indexing services, such as PubMed, as is the case for
428 some existing micro-publication communication platforms⁵⁵. Active patient
429 involvement and the ability to access, upload and delete their data enhance patient
430 autonomy and facilitate the acquisition of longitudinal patient data, further enriching
431 GMDB's repository of facial images. Similar to other natural history study data, the

432 longitudinal image and associated phenotypic meta data add significant value to the
433 understanding of disease progression in patients with facial dysmorphism⁵⁶. Moreover,
434 micro-publication encourages the recruitment of patients from underrepresented
435 populations. For example, more than 40% of all images obtained for Africans had been
436 previously unpublished. These micro-publications from unpublished images of
437 patients with underrepresented ancestries underscored the importance of GMDB
438 since they cannot be found in any medical journals.

439 The diverse ancestry data in GMDB further enabled us to investigate the
440 GestaltMatcher performance differences among ancestral groups systematically. In

441

442

443 **Table 2**, the performance disparities in the Asian and African groups were observed
444 when compared to the European group. The “Others” group showed a comparable or
445 even higher performance than the European group. The reason could be that Latin
446 Americans in the ‘Others’ group show relatively similar facial phenotypes to the
447 Europeans.

448 Our findings indicate that increasing the ancestral diversity in FAIR databases will
449 particularly benefit populations currently regarded as underprivileged. We investigated
450 how the top-1 and top-5 accuracies for the different ancestries changed when equally
451 sized groups of European or non-European patients were added to the training set.
452 Overall, the top-5 accuracy for non-European ancestral groups increased significantly
453 when the training set was expanded with non-Europeans (+11.29%). When the
454 training data were extended from only Europeans to Europeans and non-Europeans,
455 only a marginal change in the performance of the European group was observed.
456 Including more non-European patients in the gallery can also improve non-European
457 groups' performances dramatically while European performance remains roughly the
458 same (Figure 5b). The results indicate that recruiting non-European patients to support
459 the underrepresented populations is more effective than recruiting more European
460 patients, which often leads to models' extreme bias toward European ancestry.

461 The GMDB-FAIR dataset offers a transparent AI training set, which is crucial for the

462 NGP development because all FAIR data are available to the clinical and scientific
463 community. This transparency, combined with the increased representativeness of the
464 training set, helps minimise the risk of algorithmic bias, which is key for ensuring
465 respect for the fundamental right to non-discrimination⁵⁷. The high quality of the GMDB
466 data allows researchers to train, validate, and test AI in a manner that aligns with the
467 expectation in the EU AI Act and the EU Medical Device Regulation⁵⁸. Finally, the
468 controlled access and consent options as described in the Methods section not only
469 ensures respect for the fundamental right to protection of personal data⁵⁷ and EU
470 General Data Protection Regulation (GDPR)⁵⁹ compliance, but it also enabled the
471 creation of a more diverse, representative, and larger data set as people are more
472 comfortable with sharing health and genetic data, including images, under controlled
473 conditions and responsible data governance than in open access publications and
474 repositories. By this, the GMDB-FAIR dataset falls in line with other large public
475 datasets, such as ImageNet⁶⁰ for object classification or Labeled Faces in the Wild
476 (LFW)⁶¹ for face verification, which have been fundamental for deep-learning
477 technology driving computer vision over the last decade. GMDB-FAIR has been used
478 to develop many NGP approaches^{19,35–37} for predicting rare disorders after the first
479 usage in GestaltMatcher in 2022. Moreover, GMDB-FAIR data can be used in the
480 research platform (Supplementary Note) to validate the results shown in the published
481 works^{39–53} that provides transparency to the researcher using GestaltMatcher and the
482 probability to extend the existing research with the user's additional data.

483 Due to variability in facial phenotypes secondary to ancestry, diverse reference image
484 databases are crucial in order to enable clinicians to learn about the phenotypic
485 variability in facial dysmorphism within a given disorder. While efforts have been made
486 to create an atlas of human malformations that addresses the issue of ancestral
487 diversity, this remains limited to only a few disorders²⁰. With GMDB-FAIR, we created
488 a large-scale dataset that can be searched for disorders or genes of interest in the
489 GMDB gallery view (Figure 2c, Figure 3b), which provides clinicians with a
490 comprehensive selection of patient images from different ancestries at a glance,
491 thereby eliminating the need for extensive literature searches. In addition, it facilitates
492 facial phenotype comparisons within a given disorder among different ancestries
493 (Supplementary Note). GMDB also represents a valuable teaching tool for training
494 students and residents to recognize disorders based on facial features.

495 To conclude, GMDB is a medical imaging database for rare disorders that
496 encompasses diverse populations. The FAIR data will serve as reference material for
497 clinicians that facilitates learning about facial dysmorphism across ancestries, and as
498 a transparent training and benchmarking dataset for advancing the NGP approach.
499 While we show improved performance for the underrepresented populations, it is
500 important to point out that the performance is far from the optimum that can be
501 achieved by collecting more diverse data. We envision that the gap between the
502 European ancestral group and the underrepresented ancestries can be mitigated by
503 micro-publications in the future, and this will result in substantially improved support
504 for underrepresented populations.

505 **Methods**

506 **Implementation of the online GMDB platform**

507 The online platform was built using Ruby on Rails in order to allow users to input
508 images and other patient data. A database was set up using MySQL to store the
509 patient data. GMDB is hosted physically in the University Hospital of Bonn and is
510 maintained by Arbeitsgemeinschaft für Gen-Diagnostik e.V. (AGD), which is a non-
511 profit organization for genomic research. The service is funded by membership fees
512 of the AGD and donations from the Eva-Luise und Horst Köhler Foundation and the
513 Wirtgen Foundation.

514 **Image data and meta data stored in GMDB**

515 An entry in GMDB consists of a medical image such as a portrait, X-ray, or fundoscopy
516 and machine-readable meta information containing: 1) demographic data (including
517 sex, age, and ancestry); 2) the molecularly confirmed diagnosis (OMIM index⁶²); 3) the
518 disease-causing mutation reported in Human Genome Variation Society format⁶³
519 (HGVS) or International System for Human Cytogenomic Nomenclature⁶⁴ (ISCN) with
520 test method and zygosity; and 4) the clinical feature encoded in HPO terminology³⁴
521 (Figure 2b). When submitting data, clinicians are also asked to state their expert
522 opinion concerning the distinctiveness of a phenotype: They are asked to score
523 whether the medical imaging data was supportive (1), important (2), or key (3) in
524 establishing the clinical diagnosis. Computer scientists can use this information to
525 interpret the performance of their AI¹⁵.

526 **Digital consent form and patient-centered data upload**

527 To facilitate faster retrospective patient recruitment, a digital consent form has been
528 implemented, which allows patients to select conditions for storing their data within the
529 database and enables the provision of their signature online. To address the specific
530 requests of patients, this feature was further developed in close collaboration with
531 patient support groups, e.g., the German Smith-Magenis Syndrome patient
532 organization Sirius e.V. Patients can access their own cases and provide or withdraw
533 their consent online. They can also upload images themselves, which greatly simplifies
534 the curation process for longitudinal image data and other prospective data. The fact
535 that documents such as letters from clinicians or laboratory results can also be
536 uploaded, while only being visible to the responsible clinician, makes it possible to
537 obtain molecular and phenotype information on patients recruited retrospectively from
538 patient support groups. This digital consent is developed in such a way that it could
539 also, in principle, be used as a dynamic consent model in the future⁶⁵. The consent
540 form is available in German and English, and other languages will be incorporated in
541 the near future. Please find them in Supplementary Note (Digital consent, and
542 Supplementary Figures 5 and 6) for more details.

543 **Data curation**

544 The curated data can be broadly categorized as retrospective and prospective.
545 Retrospective refers primarily to data collected from the literature or from similar
546 projects with global consent for data sharing (e.g., Minerva&Me⁶⁶). For cases curated
547 from the literature, the DOI and PubMed ID as well as the contact details of the
548 corresponding author were collected in order to clarify whether reuse is possible while
549 respecting intellectual property rights. Following the provision of written informed
550 consent, our collaboration partners, clinicians from around the world (Figure 1a and
551 the co-authors), also recruited patients with an established diagnosis from within their
552 clinical practice or from patient support groups. Prospective curation refers to the
553 collection of further images or metadata over time. This can be done by the attending
554 clinician after subsequent consultations, or by the patients themselves.

555 The curation process can be broadly subdivided into three phases. First, medical
556 students in their final year annotated cases from the literature, mainly searched

557 PubMed and Google Scholar for publications with images of patients with facial
558 dysmorphism and monogenic molecular diagnosis.

559 Second, solved patients were recruited from patient support groups. Included patients
560 were allowed to upload and delete images and findings autonomously and access
561 their data at any time. To develop a patient-centered, user-friendly platform and
562 strengthen patient autonomy, feedback was obtained from the recruited patients
563 during this phase in order to determine whether any adjustments to the process were
564 required.

565 In the third phase, the database was expanded via international collaborations with
566 clinicians from different continents. Initially, this focused on patients who had already
567 been solved but had not yet been published in order to improve the AI's performance.
568 However, as we progressed, more clinicians shared their unsolved cases with the
569 scientific community. GMDB then started focusing on facial portraits of patients with
570 rare monogenic diseases, and is now dominated by, but not limited to, such cases.
571 Later in the curation process, we also annotated cytogenetic disorders with facial
572 dysmorphism. In addition to these clinicians, the medical students continued to
573 annotate data from the literature.

574 **Digital Object Identifier assignment**

575 After data submission, the respective case is immediately published on the website.
576 Subsequently, the author has the option of generating a DOI in order to create a citable
577 micro-publication⁵⁴. To do this, clinicians must, after uploading the required data and
578 metadata, enter their own personal identifier (e.g., ORCID), specify all other scientists
579 or clinicians involved in this case, and provide a title and an abstract. To ensure the
580 credibility and reliability of the published data, this process will adhere to a rigorous
581 review similar to that described by Raciti et al.⁵⁵. The DOIs are created and managed
582 by the University and State Library of Bonn using the DataCite Application
583 Programming Interface (API) (<https://datacite.org>).

584 Additionally, a dedicated landing page will be created for each case, according to the
585 specifications of the DataCite metadata schema (Supplementary Figure 2). The
586 landing page is accessible via the generated DOI, even for individuals without access
587 to GMDB or those who are not logged in. The landing page contains the full citation

588 with the DOI as a link, the abstract, and a description of the case data. No phenotypic
589 information, HPO terms, or images are available. However, the landing page indicates
590 how many images the micropublication contains.

591 **Main components of the GMDB online platform**

592 The GMDB consists of three main components that can in principle be utilized by
593 registered users (Figure 2c). 1) Search: Clinicians can use the Gallery view to search
594 the GMDB for disorders or genes of interest and get all patients matching this search
595 criterion displayed in the database at a glance. 2) Analyze: Clinicians and scientists
596 can use the GMDB-FAIR data to perform similarity comparisons of cohorts with
597 GestaltMatcher within the research platform of GMDB. 3) Train: The GMDB-FAIR
598 dataset that can be used by external researchers to train NGP tools. More detailed
599 information on these features can be found in the Supplement Note.

600 **GMDB datasets**

601 All analyses performed in this paper are based on GMDB-FAIR data (v1.1.0). But
602 actually, the GMDB consists of the GMDB-FAIR dataset and the GMDB-private set
603 (Supplementary Note and Supplementary Figures 7 and 8). We introduced this
604 distinction because it is known that patient consent to data sharing is higher when not
605 shared with a broad mass, but only for a specific study⁶⁷. However, many patients
606 agree to controlled access for the general scientific community to advance research⁶⁷.
607 For this reason, patients can decide whether they want to be part of only the GMDB-
608 private set for AI training or agree to be part of the FAIR data set.

609 The website displays the statistics to the public, showing how many patients are in the
610 database and how many disorders and disease genes have been curated. When the
611 user has the link to a specific case in the GMDB (e.g., from a publication in which the
612 original image may not be branched, but a link to the case is given in the GMDB), if
613 the user is not logged in, the landing page for the case will show how many images
614 and metadata are available for the case. Only sex and ancestry, as well as the disease
615 gene, are given. If it is a case report published with a DOI in the GMDB, the
616 corresponding title and abstract of the case can also be viewed. The remaining data
617 can only be viewed after logging in. To visualize the images, the user has to log in to
618 the platform.

619 **GMDB-FAIR data set**

620 The FAIR data set (Supplementary Figure 7b) is accessible to the scientific community.
621 Data comes from publications and from clinicians or patients themselves. However,
622 the case is accessible in the Gallery view for all registered users of the GMDB, and
623 the data sheet with all relevant data and metadata can be viewed. It is also available
624 to all users of the GMDB to perform similarity comparisons of cohorts in the research
625 platform (Supplementary Note). The data is used for the GestaltMatcher training and
626 test set but can also be made available to other scientists to train and test their AI after
627 they have applied to us with an Institutional Review Board (IRB)-approved study and
628 proposal.

629 **Data Governance and Ethical, Legal and Social Implications of GMDB**

630 Ethical approval for the GMDB was granted by the IRB of the University of Bonn, and
631 all patients have given informed written consent to participate. During the
632 GestaltMatcher consent procedure, patients can also indicate whether they agree to
633 the use of the images in presentations, teaching activities, or in publications in other
634 journals. This differentiation from other journals is important since patients/parents
635 show less willingness to consent to publication in open-access journals than to
636 publication in access-controlled databases that are not publicly accessible⁶⁷.

637 The GMDB has four different levels of data access (Supplementary Figure 8): 1) The
638 public data, which includes a summary of the GMDB statistics on the website and a
639 landing page for case reports with DOI (Supplementary Figure 2), requires no login
640 and is openly accessible. 2) The FAIR data, which can be viewed with a GMDB user
641 account, and in principle, downloaded by external AI researchers. 3) The restricted
642 data, which is not accessible to GMDB users and external AI researchers and can only
643 be used to train the GestaltMatcher AI. 4) Patient-shared data: Patients can only view
644 their own case and upload data if they are invited to do so by the attending clinician.

645 External scientist in the field of AI can apply to download of GMDB-FAIR data for the
646 development of NGP approaches. Prerequisites for this are IRB approval and
647 submission of a proposal to info@gestaltmatcher.org. In addition, external scientists
648 must sign and adhere to the GDPR. The Advisory Board will conduct a thorough review
649 of all applications. If the majority of the members of the Board approve the application,

650 access (under the extent permissible by law) will be granted to applicants within two
651 to three weeks.

652 **Advisory Board**

653 Advisory Board comprises the following co-authors: Benjamin D. Solomon, Koen
654 Devriendt, Shahida Moosa, Christian Netzer, Martin Mücke, Christian Schaaf, Alain
655 Verloes, Christoffer Nellåker, Markus M. Nöthen, Gholson J. Lyon, Aleksandra Jezela-
656 Stanek, and Karen W. Gripp.

657 **HPO-defined symptom groups**

658 In one of our previous works⁶⁸, twelve distinct and non-overlapping categories of HPO
659 terms were defined by clinical experts (“HPO defined symptom groups”). All GMDB
660 cases for which HPO terms were annotated were then assigned to each of those
661 groups, if at least one of the HPO terms in this group was annotated; i.e., each GMDB
662 case can be assigned to several HPO-defined symptom groups. For each case, the
663 most pronounced HPO-defined symptom group was defined as the single group
664 comprising the largest number of the case’s annotated HPO terms. The HPO-defined
665 symptom group “Others” was only assigned as the leading HPO-defined symptom
666 group if no other HPO-defined symptom group was present for the case.

667 Phenotypic similarity between cases was calculated using the R-package
668 ontologySimilarity (version 2.5). Pairwise similarities were calculated for the combined
669 data set of GMDB cases with HPO terms (n=4,474), the TRANSLATE-NAMSE exome
670 sequencing data set (n=1,577), and data on known diseases and their clinical features
671 downloaded from the HPO website (n=7,765,
672 <https://hpo.jax.org/app/download/annotation>, file: genes_to_phenotype.txt,
673 downloaded on 10 April 2021). The resulting distance matrix was projected in a four-
674 dimensional space via Uniform Manifold Approximation and Projection (UMAP). The
675 first two dimensions were plotted using ggplot2 (version 3.4.4). To analyze which
676 HPO-defined symptom groups occur jointly, the proportion of patients assigned to the
677 first group that were also assigned to the second group was assessed. All analyses
678 were conducted in R (version 4.3.2).

679 **GestaltMatcher Algorithm**

680 DeepGestalt¹⁷ is a deep learning-based NGP tool using frontal face photos to classify
681 up to 216 syndromes it has seen during training. However, it needed a lot of training
682 data to achieve a reasonable performance on these syndromes. That also meant it
683 could not classify unseen syndromes during training (ultra-rare syndromes). This led
684 to the development of GestaltMatcher¹⁵, which uses a clustering approach. As such,
685 if at least one image of the sought-after syndrome is in the gallery set, a test image
686 can be matched to/clustered with that image using some similarity metric. Later this
687 approach was further enhanced by Hustinx et al.¹⁹, using a more recent architecture
688 (iResNet) and training loss (ArcFace Loss), as well as test-time augmentation and a
689 model ensemble to improve robustness. That is also the approach we used for our
690 experiments. Thus, for fine-tuning we utilized the Adam optimizer, cross-entropy loss,
691 and class weighting to deal with the imbalance in data availability between disorders.

692 The overall idea behind the methodology is to train a classifier on a more frequent
693 subset of the syndromes, achieving a model that generalizes well on those seen
694 syndromes. In practice, the authors of both papers decided to use syndromes with at
695 least seven patients as the training set for this classifier. Thereafter, everything up to
696 the penultimate layer of the classifier is used as an encoder, obtaining feature
697 embeddings of images of interest. These could be images for the gallery set or images
698 for the test set.

699 The aforementioned gallery set is the set of images (and their feature embeddings)
700 with known syndromes. This can include the syndromes used for training (seen) and
701 syndromes with too few images to train on (unseen). The theory is that similar facial
702 phenotypes form clusters in the feature space, which is spanned by the feature
703 embeddings in 512 dimensions and which we refer to as clinical face phenotype space.
704 The similarity between images and clusters is computed using the cosine distance,
705 where a lower distance implies a higher similarity. Contrary to the approach by
706 Gurovich et al.¹⁷, this approach can easily increase support for ultra-rare syndromes.
707 The quality and diversity of the gallery set is crucial for this approach to match test
708 images to clusters in the gallery set.

709 **Performance metric (top-k accuracy)**

710 The applied performance metric was top-k accuracy. Top-1 indicates that the disorder
711 was correctly classified as the first guess, while top-5 indicates the correct class was
712 in the first five guesses. We reported top-k accuracies (k=1, 5, 10, and 30) as the
713 performance readout.

714 **Ancestry analysis**

715 The genetic ancestry of each individual was documented as precisely as possible
716 using self-reported data. For instance, if an individual was born in Germany and all of
717 the respective grandparents also originated from there, this individual was assigned
718 to Germany (country) and Europe (continent). The same approach was used for all
719 individuals with no self-reported migration history in previous generations. For
720 individuals with mixed ancestry, the respective ancestries were combined. For
721 example, an individual with a father from Gambia and a mother from Eastern Europe
722 was assigned European-African mixed ancestry.

723 The performance of GestaltMatcher is highly dependent on the training set and the
724 gallery set. To investigate the impact of incorporating diverse ancestry on the
725 performance, we have therefore conducted two sets of experiments for those two
726 components, respectively. First, we analyzed the influence on the models'
727 performance when including only European versus both European and non-European
728 data into the training set. And second, we analyzed the same performance when
729 iteratively increasing the amount of non-European data into the gallery set.

730 In the first experiment, a subset of images of European patients (EU) was extended
731 by either the inclusion of a different subset of images of European patients (EU*), or a
732 subset of patients with non-European ancestries (non-EU) (Supplementary Figure 9).
733 Random sampling of these subsets was performed five times. EU consisted of on
734 average 3,139.2 images, and EU* comprised on average 1,567.6 images. First, the
735 model was trained on the EU + EU* set containing on average 4,706.8 images of
736 patients of solely European ancestry. For EU + non-EU, a subset containing on
737 average 1,567.6 images of patients with any non-European ancestry was used,
738 totaling to 4,706.8 images. The experiment design ensured the maintenance of the
739 same distribution of disorders as that found in the training data.

740 The model was fine-tuned for 50 epochs on subsets EU + EU* and EU + non-EU of
741 GMDB (v1.1.0). All other hyperparameters were left unchanged. It is important to note
742 that the model was not tasked with learning to classify the ancestry, only with learning
743 to classify the disorder.

744 Post-training, the models' performances were measured on the same evaluation set,
745 containing images of patients with diverse ancestral backgrounds. This evaluation set
746 consisted of 649 images and was sampled in such a manner that there was no overlap
747 between patients or images in any subset. Top-k accuracy was averaged over each
748 ancestry rather than each image in order to address the imbalance in ancestry
749 frequency. As such, the performance of any infrequent group weighed equally with
750 those of the more frequent groups.

751 In the second set of experiments, we trained the models of Hustinx et al.¹⁹ using the
752 GMDB-FAIR training set, including different proportions of non-EU data for the gallery
753 set. We compared the performance of the syndromes our models have seen during
754 training. For completeness, Table 1 shows the top-k accuracy (over all images) for
755 different categories (sex, ancestry, and age range) using the entire gallery set (100%
756 EU + 100% non-EU). For the experiments, we computed the performance when
757 including different proportions of non-EU data, extending the gallery set by +10% per
758 iteration. This experiment was repeated tenfold, randomly sampling patients with
759 different ancestries and all their photos for the gallery set. As such, at 0%, we include
760 only data from EU patients in the gallery set, and at 100%, we include all patient data
761 for the relevant syndromes.

762 We further computed the performance on syndromes that occur in both the European-
763 group and each non-European group to more accurately reflect the performance
764 differences, avoiding the imbalance between offered support for each ancestral group.

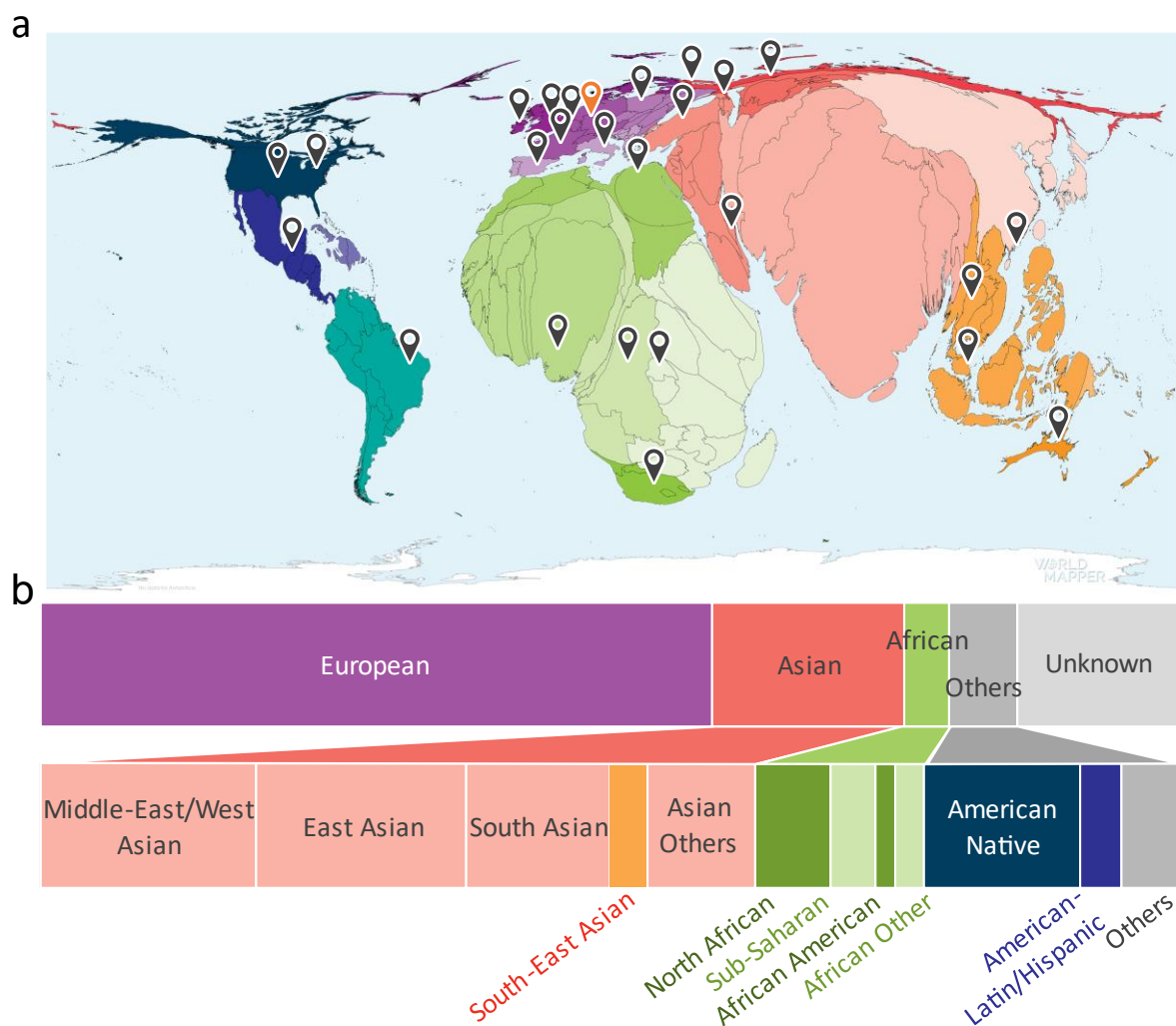
765 Data and code availability

766 GMDB-FAIR can be downloaded in GMDB after the application is approved by the
767 advisory board. Please find more details in the Data Governance and ELSI section.
768 Code is available in the GitHub repository ([github.com/igsb/GestaltMatcher-
769 Arc/tree/gmdb](https://github.com/igsb/GestaltMatcher-Arc/tree/gmdb)).

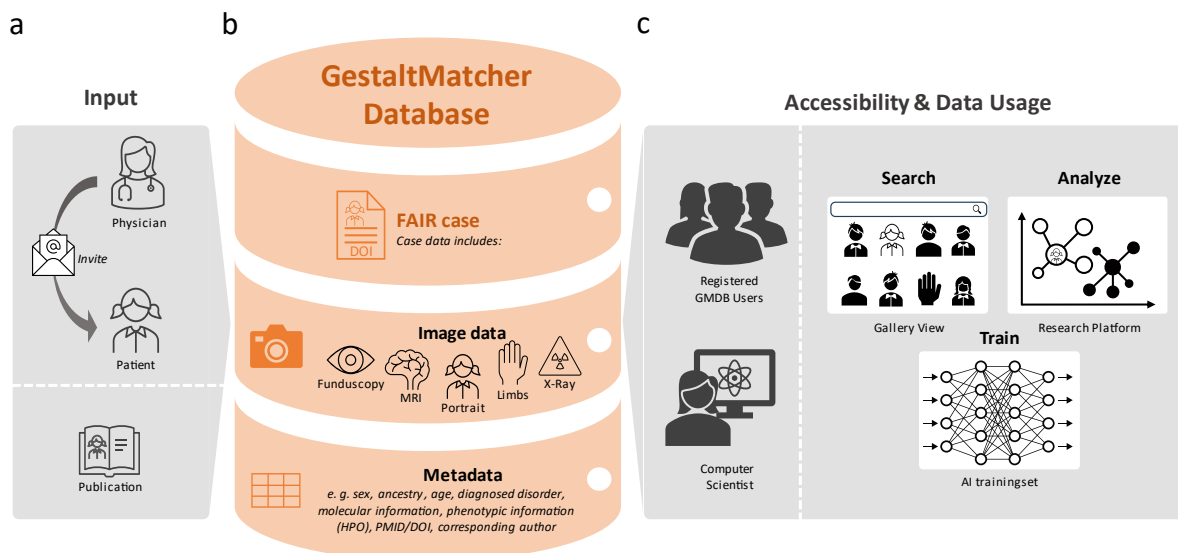
770 Acknowledgments

771 This research was supported in part by the Intramural Research Program of the
772 National Human Genome Research Institute, National Institutes of Health, the United
773 States of America. Sofia Douzgou Houge was supported by the Norwegian National
774 Advisory Unit on Rare Disorders (grant number #43066). Tahsin Stefan Barakat was
775 supported by the Netherlands Organisation for Scientific Research (ZonMw Vidi, grant
776 09150172110002). Heidi Beate Bentzen is supported by EU grant 101071203 and
777 Research Council of Norway grants 322672 and 324278. Nina-Maria Wilpert was
778 supported by the DFG Research Unit 2841 “Beyond the Exome” and is a participant
779 in the BIH Charité Junior Clinician Scientist Program funded by the Charité -
780 Universitätsmedizin Berlin, the Berlin Institute of Health at Charité (BIH), the
781 Alliance4Rare, and the Berliner Sparkassenstiftung Medizin. Cristina Dias was
782 supported by the Wellcome Trust [grant number 209568/Z/17/Z]. The authors thank
783 the Asia Pacific Society of Human Genetics, the Wirtgen Foundation, the Eva Luise
784 und Horst Köhler Foundation, Kabuki Syndrome Foundation, Kleefstra support group,
785 German Smith-Magenis Syndrome patient organization Sirius e.V. and the Focus
786 Foundation for their support.

787 Figures



788
 789 **Figure 1: a)** Birth rate distribution worldwide. The size of country is scaled in
 790 accordance with the respective birth rate. The map indicates countries from which
 791 unpublished images were obtained (source: <https://worldmapper.org/faq/>, modified).
 792 **b)** Distribution of ancestry groups in GestaltMatcher Database. 16% of the patients
 793 without ancestral information were categorized as Unknown. The breakdown of
 794 ancestries in the dataset with known ancestry is as follows: European 67%, Asian 19%,
 795 African 7%, and Others 7%.



796

797 **Figure 2: GestaltMatcher Database (GMDB) Architecture and Dataflow. a)**

798 Retrospective data are collected from the literature and annotated by data curators or
799 are uploaded by collaborating attending clinician. Patients can also upload images of

800 their own cases, incorporate prospective data, and view their own data at any time. **b)**

801 The data (multimodal image data, including portrait images as well as magnetic
802 resonance imaging, X-ray, funduscopy and extremity images) are stored in the GMDB

803 (MySQL database) together with the relevant meta information (such as sex, age,
804 ancestry, molecular, and phenotypic information). **c)** Registered users can view and

805 search the FAIR data in the GMDB Gallery. The patient image can also be analyzed
806 using the Next-Generation Phenotyping tool GestaltMatcher within the Research

807 Platform. In addition, once their application has been approved by the Advisory Board,
808 external computer scientists can use the GMDB-FAIR data set for training purposes

809 for their projects.

810

811 **Figure 3: An example case presentation of a FAIR case with a Digital Object**

812 **Identifier (DOI).** **a)** FAIR cases in the GestaltMatcher Database (GMDB) are displayed
813 to GMDB users via the data sheet. Each FAIR case can also be assigned a DOI in

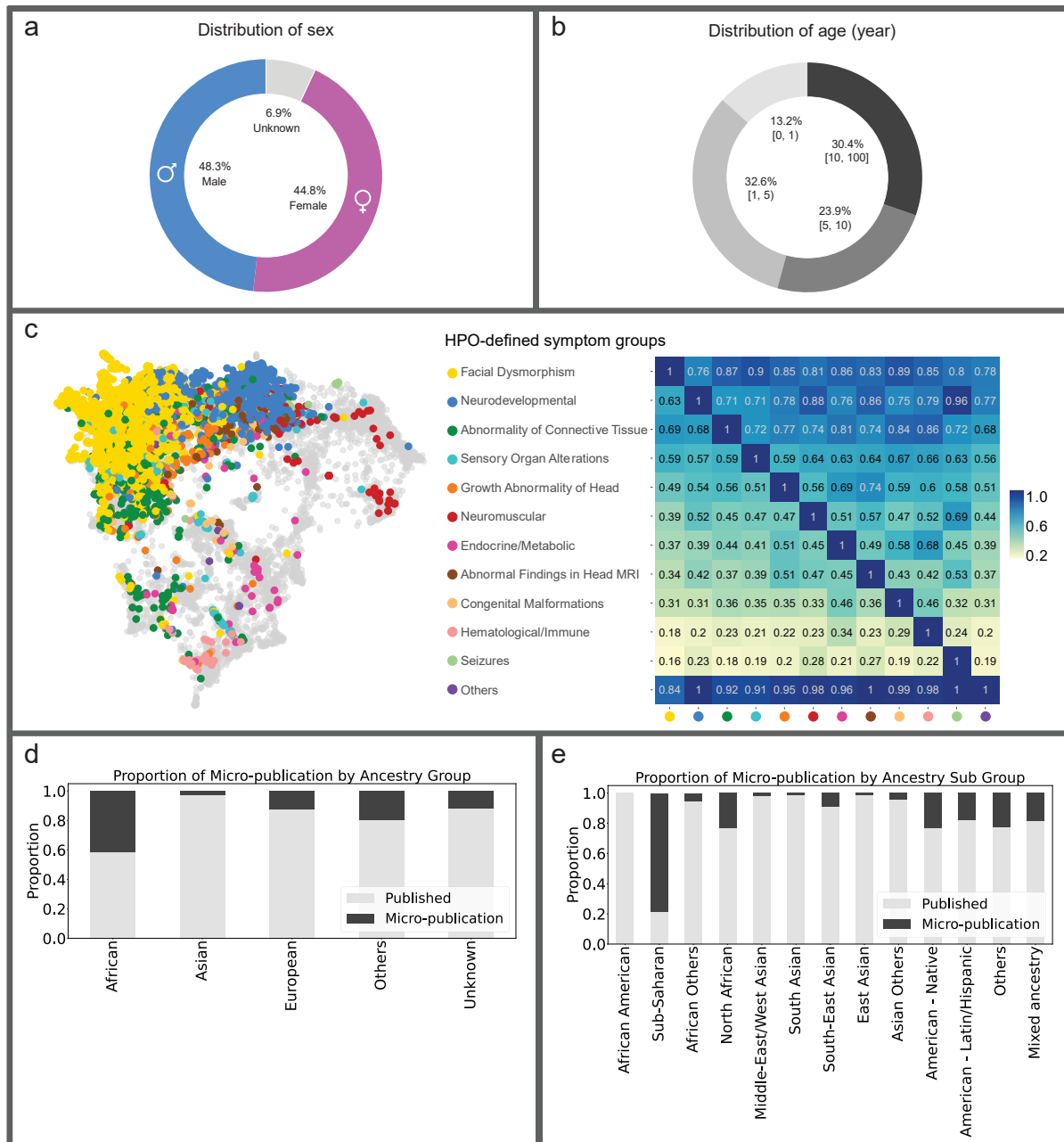
814 order to render it a citable micro-publication. This micro-publication contains the image
815 data and metadata, including demographic, molecular, and phenotype information.

816 The dynamic nature of the GMDB case report enables longitudinal image data storage
817 even after initial publication, which is not possible in conventional journals. **b)** After

818 uploading, case reports can be viewed and searched by other users in the Gallery

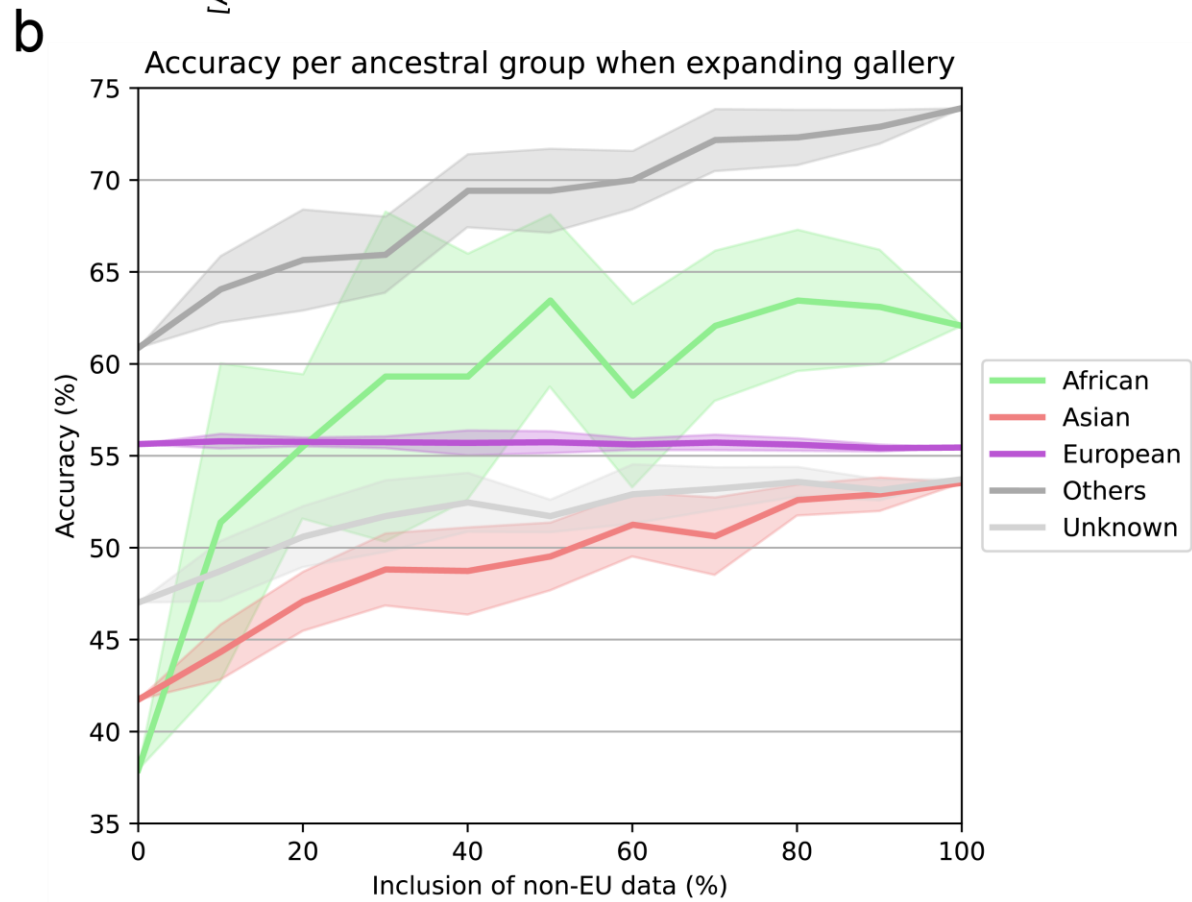
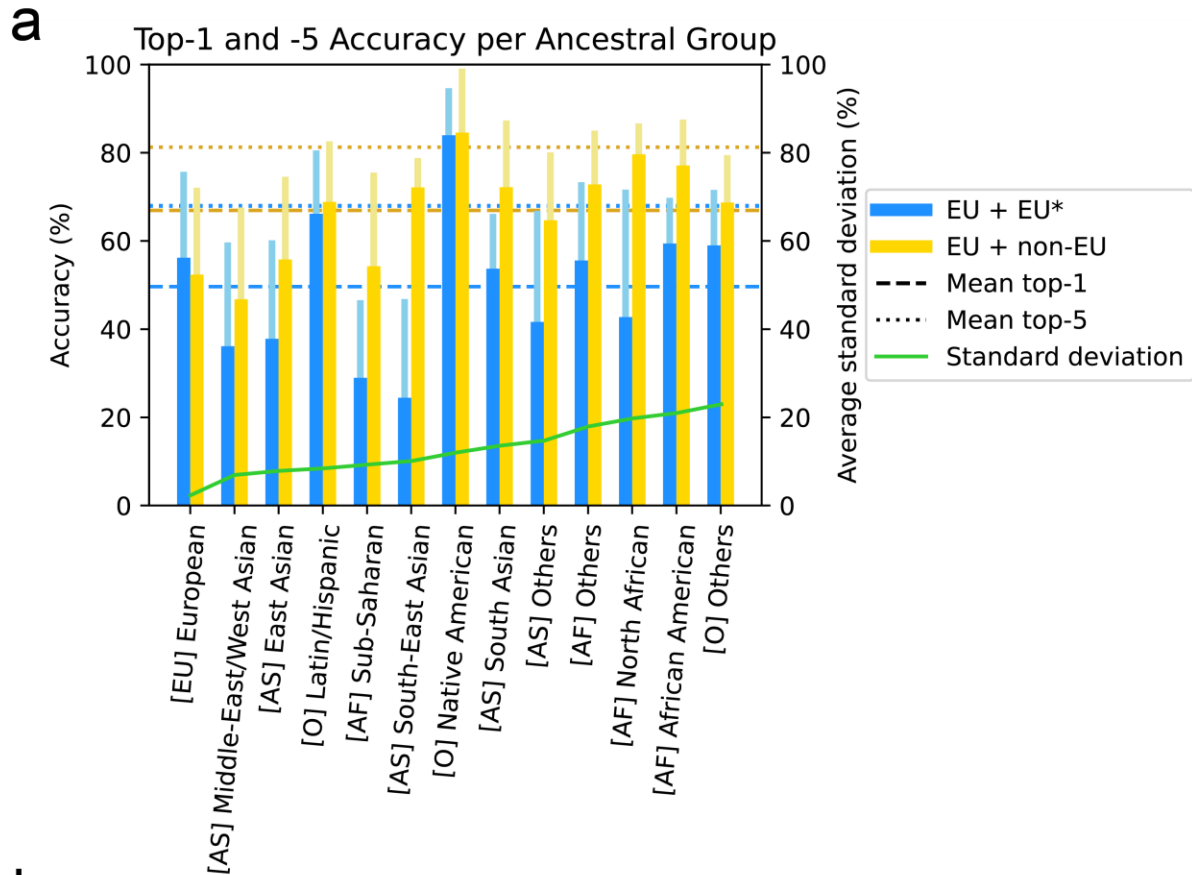
819 view. **c)** The image data can also be used for inter-cohort comparisons of the gestalt
 820 scores within the research platform.

821



822 **Figure 4: Overview of the GestaltMatcher Database (GMDB)-FAIR dataset. a)** Sex
 823 distribution. Number of images shown in brackets. **b)** Distribution of patient age in
 824 years. **c)** Left: Two-dimensional representation of phenotypic similarities between
 825 patients, as calculated on the basis of Human Phenotype Ontology (HPO) terms via
 826 Uniform Manifold Approximation and Projection (UMAP). HPO terms were annotated
 827 for 4,474 individuals in the GMDB, and expert clinicians defined twelve distinct HPO-
 828

829 defined symptom groups. Based on the annotated HPO terms, each case was
830 assigned to one or more HPO-defined symptom groups. All OMIM diseases were
831 included, using their HPO annotations (gray background dots) as a reference. GMDB
832 cases are color-coded according to their most pronounced HPO-defined symptom
833 group, i.e., the group that includes the majority of their HPO terms. The dataset is
834 dominated by two major clusters (facial dysmorphism in yellow and
835 neurodevelopmental in blue) but shows cases from across the complete disease
836 landscape. Right: Heatmap of the proportion of GMDB individuals within the HPO-
837 defined symptom group on the X-axis who are also assigned to the HPO-defined
838 symptom group on the Y-axis. Notably, facial dysmorphism is present in at least 70%
839 of the cases of each HPO-defined symptom group. **d)** Proportion of the unpublished
840 and published images in each ancestry group. **e)** Proportion of the unpublished and
841 published images in each sub-ancestry group.



842

843 **Figure 5: Performance of ancestry analysis. a)** Top-1 and top-5 accuracy of
844 GestaltMatchers' disorder classification accuracy per ancestral group. Top-1 and top-
845 5 accuracy of the models' disorder classification accuracy per ancestral group, where
846 (blue) belongs to the EU only subset, and (yellow) belongs to the diverse subset. Each
847 wide, darker bar and each light, thinner bar indicate the top-1 and top-5 accuracy per
848 ancestral group, respectively. The horizontal dashed lines and dotted lines indicate
849 the top-1 and top-5 overall accuracy averaged over all ancestral groups, respectively.
850 The order of the ancestry group in the x-axis is ranked according to standard deviation
851 between top-1 accuracies of the 5-fold experiment. **b)** Top-1 accuracy of
852 GestaltMatcher when including different proportion of non-European patients in the
853 gallery. The x-axis is the proportion of non-European data included in the gallery. The
854 y-axis is the top-1 accuracy. The colored region along the line indicates the standard
855 deviation.

856 Tables

857 **Table 1: Performance of GestaltMatcher on different categories of sex, ancestry,**
 858 **and age.** The top-1, top-5, top-10, and top-30 accuracy are reported. For the top-1 to
 859 top-30 columns, the best performance in each category is boldfaced. In the ancestry
 860 category, the sampling influences European and other ancestry groups' performance
 861 due to the significant difference in the test image size. They may evaluate the different
 862 sets of disorders. We, therefore, presented the performance of the overlapped
 863 disorders in

864
865

866 **Table 2.** In the age category, the notation $[x, y)$ represents a half-open interval, which
 867 includes the starting point x but excludes the endpoint y . For example, $[0, 1)$ years
 868 range from birth but do not include one year old.

Category	Test images	Top-1	Top-5	Top-10	Top-30	
Overall	882	56.58%	76.08%	82.61%	90.36%	
Ancestry	African	29	62.07%	82.76%	82.76%	86.21%
	Asian	127	53.54%	78.74%	85.04%	89.76%
	European	523	55.45%	75.14%	82.60%	90.25%
	Others	69	73.91%	81.16%	81.16%	92.75%
	Unknown	134	53.73%	73.13%	81.34%	91.04%
Sex	Male	419	55.37%	74.22%	80.67%	88.78%
	Female	393	55.98%	75.83%	83.21%	91.09%
	Unknown	70	67.14%	88.57%	91.43%	95.71%
Age	$[0, 1)$ years	53	52.83%	71.70%	79.25%	90.57%
	$[1, 5)$ years	137	56.20%	75.91%	81.02%	90.51%
	$[5, 10)$ years	115	57.39%	83.48%	86.09%	90.43%
	$[10, \infty)$ years	165	58.18%	71.51%	77.58%	85.45%
	Unknown	412	56.31%	76.46%	84.71%	92.23%

869
870
871

872 **Table 2: Performance comparison between European and other ancestry groups**
 873 **on the overlapping disorders.** This table is an extension of the ancestry section in
 874 Table 1, taking the overlapped disorders between European and other ancestry
 875 groups. Each category compares European and non-European ancestry groups'
 876 performance on the same set of disorders. The number of overlapped disorders is
 877 reported in the 'Disorders' column. In comparing African and European groups, six
 878 disorders exist in the test sets of both ancestry groups. The top-1, top-5, top-10, and
 879 top-30 accuracy are reported. For the top-1 to top-30 columns, the best performance
 880 in each category is boldfaced.

Category	Disorders	Test images	Top-1	Top-5	Top-10	Top-30
[African, European]	African	14	64.29%	85.71%	85.71%	92.86%
	European	32	81.25%	96.88%	96.88%	100.00%
[Asian, European]	Asian	83	57.83%	79.52%	84.34%	87.95%
	European	139	64.75%	82.73%	88.49%	91.37%
[Others, European]	Others	53	81.13%	90.57%	90.57%	100.00%
	European	115	69.56%	84.35%	93.91%	96.52%
[Unknown, European]	Unknown	77	59.74%	81.81%	88.31%	96.10%
	European	170	62.35%	81.18%	89.41%	94.12%

881

882

883

884 **Table 3: Training accuracy with EU + non-EU and EU + EU* datasets.** Within the
 885 European training row, numbers annotated with * in brackets indicate the training
 886 images from EU + EU. Higher top-1 and top-5 accuracies between EU + EU* and EU
 887 + non-EU training are denoted in bold.

	Number of images		Performance EU + non-EU		Performance EU + EU*	
	Training	Testing	Top-1	Top-5	Top-1	Top-5
European	(4706.2 ± 24.4)* 3139.2 ± 15.1	444.6 ± 22.2	52.35 ± 2.30%	72.05 ± 2.66%	56.17 ± 2.27%	75.66 ± 2.70%
East Asian	283.2 ± 5.0	31 ± 6.2	55.78 ± 10.25%	74.56 ± 5.90%	37.77 ± 5.45%	60.13 ± 5.77%
Latin/Hispanic	257.8 ± 7.0	28.4 ± 4.7	68.86 ± 8.92%	82.56 ± 7.77%	66.16 ± 7.89%	80.51 ± 6.58%
Middle-East/ West Asian	211.2 ± 6.8	30 ± 5.8	46.76 ± 7.01%	67.59 ± 7.52%	36.10 ± 6.81%	59.67 ± 3.68%

South Asian	200.2 ± 5.4	18.8 ± 2.7	72.15 ± 12.24%	87.32 ± 10.04%	53.70 ± 14.86%	66.13 ± 13.40%
Asian Others	170.6 ± 2.4	16.4 ± 4.1	64.66 ± 10.93%	80.06 ± 13.87%	41.64 ± 18.51%	66.84 ± 11.87%
Sub-Saharan	119 ± 2.7	18.2 ± 3.4	54.23 ± 7.32%	75.49 ± 15.27%	28.91 ± 11.18%	46.55 ± 11.71%
North African	64.8 ± 3.2	7.4 ± 1.6	79.64 ± 20.19%	86.64 ± 14.62%	42.71 ± 19.34%	71.64 ± 18.38%
Native American	63.2 ± 6.8	14.8 ± 1.6	84.55 ± 12.77%	99.09 ± 1.82%	83.94 ± 11.18%	94.65 ± 8.62%
African Others	53.2 ± 2.0	6.2 ± 2.2	72.78 ± 18.29%	85.00 ± 13.33%	55.56 ± 17.57%	73.33 ± 22.61%
South-East Asian	51.4 ± 2.0	5.4 ± 1.3	72.12 ± 13.36%	78.81 ± 24.61%	24.40 ± 6.76%	46.83 ± 17.97%
Others	54.6 ± 2.3	6 ± 2.9	68.71 ± 23.67%	79.43 ± 22.38%	59.00 ± 22.35%	71.57 ± 21.09%
African American	38.4 ± 4.3	3.8 ± 2.6	77.08 ± 18.04%	87.50 ± 21.65%	59.38 ± 24.00%	69.79 ± 18.49%
Overall	4706.8 ± 26.7	631 ± 23.8	66.90%	81.24%	49.65%	67.95%

888
889
890

891 References

- 892 1. Hart, T. C. & Hart, P. S. Genetic studies of craniofacial anomalies: clinical implications
893 and applications. *Orthod. Craniofac. Res.* **12**, 212–220 (2009).
- 894 2. Lesmann, H., Klinkhammer, H. & Dr. med. Dipl. Phys. Peter M. Krawitz. The future role
895 of facial image analysis in ACMG classification guidelines. *Med. Genet.* **35**, 115–121
896 (2023).
- 897 3. Tekendo-Ngongang, C. *et al.* Rubinstein-Taybi syndrome in diverse populations. *Am. J.*
898 *Med. Genet. A* **182**, 2939–2950 (2020).
- 899 4. Kruszka, P., Tekendo-Ngongang, C. & Muenke, M. Diversity and dysmorphology. *Curr.*
900 *Opin. Pediatr.* **31**, 702–707 (2019).
- 901 5. Hadj-Rabia, S. *et al.* Automatic recognition of the XLHED phenotype from facial images.
902 *Am. J. Med. Genet. A* **173**, 2408–2414 (2017).
- 903 6. Martínez-Abadías, N. *et al.* Facial biomarkers detect gender-specific traits for bipolar
904 disorder. *FASEB J.* **35**, (2021).
- 905 7. Fang, F., Clapham, P. J. & Chung, K. C. A systematic review of interethnic variability in
906 facial dimensions. *Plast. Reconstr. Surg.* **127**, 874–881 (2011).
- 907 8. Vorravanpreecha, N., Lertboonnum, T., Rodjanadit, R., Sriplienchan, P. & Rojnueangnit,

- 908 K. Studying Down syndrome recognition probabilities in Thai children with de-identified
909 computer-aided facial analysis. *Am. J. Med. Genet. A* **176**, 1935–1940 (2018).
- 910 9. Kruszka, P. *et al.* Down syndrome in diverse populations. *Am. J. Med. Genet. A* **173**,
911 42–53 (2017).
- 912 10. Porras, A. R., Summar, M. & Linguraru, M. G. Objective differential diagnosis of Noonan
913 and Williams-Beuren syndromes in diverse populations using quantitative facial
914 phenotyping. *Mol Genet Genomic Med* **9**, e1636 (2021).
- 915 11. Lumaka, A. *et al.* Facial dysmorphism is influenced by ethnic background of the patient
916 and of the evaluator. *Clin. Genet.* **92**, 166–171 (2017).
- 917 12. Burchard, E. G. *et al.* The importance of race and ethnic background in biomedical
918 research and clinical practice. *N. Engl. J. Med.* **348**, 1170–1175 (2003).
- 919 13. Martínez-Abadías, N. *et al.* Phenotypic evolution of human craniofacial morphology after
920 admixture: a geometric morphometrics approach. *Am. J. Phys. Anthropol.* **129**, 387–398
921 (2006).
- 922 14. Fatumo, S. *et al.* A roadmap to increase diversity in genomic studies. *Nat. Med.* **28**,
923 243–250 (2022).
- 924 15. Hsieh, T.-C. *et al.* GestaltMatcher facilitates rare disease matching using facial
925 phenotype descriptors. *Nat. Genet.* **54**, 349–357 (2022).
- 926 16. Dudding-Byth, T. *et al.* Computer face-matching technology using two-dimensional
927 photographs accurately matches the facial gestalt of unrelated individuals with the same
928 syndromic form of intellectual disability. *BMC Biotechnol.* **17**, 90 (2017).
- 929 17. Gurovich, Y. *et al.* Identifying facial phenotypes of genetic disorders using deep
930 learning. *Nat. Med.* **25**, 60–64 (2019).
- 931 18. Porras, A. R., Rosenbaum, K., Tor-Diez, C., Summar, M. & Linguraru, M. G.
932 Development and evaluation of a machine learning-based point-of-care screening tool
933 for genetic syndromes in children: a multinational retrospective study. *Lancet Digit*
934 *Health* (2021) doi:10.1016/S2589-7500(21)00137-0.
- 935 19. Hustinx, A. *et al.* Improving Deep Facial Phenotyping for Ultra-rare Disorder Verification

- 936 Using Model Ensembles. in *2023 IEEE/CVF Winter Conference on Applications of*
937 *Computer Vision (WACV)* 5007–5017 (IEEE, 2023).
- 938 20. Muenke, M., Adeyemo, A. & Kruszka, P. An electronic atlas of human malformation
939 syndromes in diverse populations. *Genet. Med.* **18**, 1085–1087 (2016).
- 940 21. Mishima, H. *et al.* Evaluation of Face2Gene using facial images of patients with
941 congenital dysmorphic syndromes recruited in Japan. *J. Hum. Genet.* **64**, 789–794
942 (2019).
- 943 22. Narayanan, D. L. *et al.* Computer-aided Facial Analysis in Diagnosing Dysmorphic
944 Syndromes in Indian Children. *Indian Pediatr.* **56**, 1017–1019 (2019).
- 945 23. Elmas, M. & Gogus, B. Success of Face Analysis Technology in Rare Genetic Diseases
946 Diagnosed by Whole-Exome Sequencing: A Single-Center Experience. *Mol. Syndromol.*
947 **11**, 4–14 (2020).
- 948 24. Hennocq, Q. *et al.* Next generation phenotyping for diagnosis and phenotype-genotype
949 correlations in Kabuki syndrome. *Sci. Rep.* **14**, 2330 (2024).
- 950 25. 1000 Genomes Project Consortium *et al.* A global reference for human genetic
951 variation. *Nature* **526**, 68–74 (2015).
- 952 26. Winter, R. M. & Baraitser, M. The London Dysmorphology Database. *J. Med. Genet.* **24**,
953 509–510 (1987).
- 954 27. Murdoch Children’s Research Institute. POSSUMweb. *POSSUMweb*
955 <https://www.possum.net.au/>.
- 956 28. Patrinos, G. P. Chapter 6 - Incentives for Human Genome Variation Data Sharing. in
957 *Human Genome Informatics* (eds. Lambert, C. G., Baker, D. J. & Patrinos, G. P.) 109–
958 129 (Academic Press, 2018).
- 959 29. Mons, B. *et al.* The value of data. *Nat. Genet.* **43**, 281–283 (2011).
- 960 30. Patrinos, G. P. *et al.* Microattribution and nanopublication as means to incentivize the
961 placement of human genome variation data into the public domain. *Hum. Mutat.* **33**,
962 1503–1512 (2012).
- 963 31. Giardine, B. *et al.* Systematic documentation and analysis of human genetic variation in

- 964 hemoglobinopathies using the microattribution approach. *Nat. Genet.* **43**, 295–301
965 (2011).
- 966 32. Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and
967 stewardship. *Sci Data* **3**, 160018 (2016).
- 968 33. Lesmann, H. & Weiland, H. Atypical presentation of a case with Noonan syndrome with
969 multiple lentiginos (Version 1). (2024) doi:10.60723/10693.
- 970 34. Robinson, P. N. *et al.* The Human Phenotype Ontology: a tool for annotating and
971 analyzing human hereditary disease. *Am. J. Hum. Genet.* **83**, 610–615 (2008).
- 972 35. Sümer, Ö., Hellmann, F., Hustinx, A., Hsieh, T.-C. & Krawitz, P. Few-Shot Meta-
973 Learning for Recognizing Facial Phenotypes of Genetic Disorders. in *Caring is Sharing*
974 – *Exploiting the Value in Data for Health and Innovation* 932–936 (IOS Press, 2023).
- 975 36. Campbell, J., Dawson, M., Zisserman, A., Xie, W. & Nellåker, C. Deep Facial
976 Phenotyping with Mixup Augmentation. in *Medical Image Understanding and Analysis*
977 133–144 (Springer Nature Switzerland, 2024).
- 978 37. Wu, D. *et al.* Multimodal Machine Learning Combining Facial Images and Clinical Texts
979 Improves Diagnosis of Rare Genetic Diseases. *arXiv [q-bio.QM]* (2023).
- 980 38. Hsieh, T.-C., Lesmann, H. & Krawitz, P. M. Facilitating the Molecular Diagnosis of Rare
981 Genetic Disorders Through Facial Phenotypic Scores. *Curr Protoc* **3**, e906 (2023).
- 982 39. Ebstein, F. *et al.* PSMC3 proteasome subunit variants are associated with
983 neurodevelopmental delay and type I interferon production. *Sci. Transl. Med.* **15**,
984 eabo3189 (2023).
- 985 40. Asif, M. *et al.* De novo variants of CSNK2B cause a new intellectual disability-
986 craniodigital syndrome by disrupting the canonical Wnt signaling pathway. *HGG Adv* **3**,
987 100111 (2022).
- 988 41. Kampmeier, A. *et al.* PHIP-associated Chung-Jansen syndrome: Report of 23 new
989 individuals. *Front Cell Dev Biol* **10**, 1020609 (2022).
- 990 42. Lyon, G. J. *et al.* Expanding the phenotypic spectrum of NAA10-related
991 neurodevelopmental syndrome and NAA15-related neurodevelopmental syndrome. *Eur.*

- 992 *J. Hum. Genet.* **31**, 824–833 (2023).
- 993 43. Aerden, M. *et al.* The neurodevelopmental and facial phenotype in individuals with a
994 TRIP12 variant. *Eur. J. Hum. Genet.* **31**, 461–468 (2023).
- 995 44. Blackburn, P. R. *et al.* Loss-of-function variants in CUL3 cause a syndromic
996 neurodevelopmental disorder. *medRxiv* (2023) doi:10.1101/2023.06.13.23290941.
- 997 45. Oppermann, H. *et al.* CUX1-related neurodevelopmental disorder: deep insights into
998 phenotype-genotype spectrum and underlying pathology. *Eur. J. Hum. Genet.* **31**,
999 1251–1260 (2023).
- 1000 46. Blackburn, P. R. *et al.* Loss-of-function variants in *CUL3* cause a syndromic
1001 neurodevelopmental disorder. *medRxiv* (2023) doi:10.1101/2023.06.13.23290941.
- 1002 47. Averdunk, L. *et al.* Biallelic variants in *CRIP1* cause a Rothmund-Thomson-like
1003 syndrome with increased cellular senescence. *Genet. Med.* **25**, 100836 (2023).
- 1004 48. Oppermann, H. *et al.* CUX1-related neurodevelopmental disorder: deep insights into
1005 phenotype-genotype spectrum and underlying pathology. *Eur. J. Hum. Genet.* (2023)
1006 doi:10.1038/s41431-023-01445-2.
- 1007 49. Schmetz, A. *et al.* Delineation of the adult phenotype of Coffin-Siris syndrome in 35
1008 individuals. *Hum. Genet.* **143**, 71–84 (2024).
- 1009 50. Küry, S. *et al.* Unveiling the crucial neuronal role of the proteasomal ATPase subunit
1010 gene *PSMC5* in neurodevelopmental proteasomopathies. *medRxiv* (2024)
1011 doi:10.1101/2024.01.13.24301174.
- 1012 51. Li, D. *et al.* Spliceosome malfunction causes neurodevelopmental disorders with
1013 overlapping features. *J. Clin. Invest.* **134**, (2024).
- 1014 52. Rigter, P. M. F. *et al.* Role of *CAMK2D* in neurodevelopment and associated conditions.
1015 *Am. J. Hum. Genet.* **111**, 364–382 (2024).
- 1016 53. Laugwitz, L. *et al.* *ZSCAN10* deficiency causes a neurodevelopmental disorder with
1017 characteristic oto-facial malformations. *Brain* (2024) doi:10.1093/brain/awae058.
- 1018 54. Clark, T., Ciccarese, P. N. & Goble, C. A. Micropublications: a semantic model for
1019 claims, evidence, arguments and annotations in biomedical communications. *J. Biomed.*

- 1020 *Semantics* **5**, 28 (2014).
- 1021 55. Raciti, D., Yook, K., Harris, T. W., Schedl, T. & Sternberg, P. W. Micropublication:
1022 incentivizing community curation and placing unpublished data into the public domain.
1023 *Database* **2018**, (2018).
- 1024 56. Liu, J. *et al.* Natural History and Real-World Data in Rare Diseases: Applications,
1025 Limitations, and Future Perspectives. *J. Clin. Pharmacol.* **62 Suppl 2**, S38–S55 (2022).
- 1026 57. European Union. Charter of Fundamental Rights of the European Union, 2016. EUR-
1027 Lex. <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=celex%3A12016P%2FTXT>
1028 (2016).
- 1029 58. Regulation (EU) 2017/745 of the European Parliament and of the Council of 5 April
1030 2017 on medical devices, amending Directive 2001/83/EC, Regulation (EC) No
1031 178/2002 and Regulation (EC) No 1223/2009 and repealing Council Directives
1032 90/385/EEC and 93/42/EEC. [https://eur-lex.europa.eu/legal-](https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX%3A32017R0745)
1033 [content/EN/TXT/?uri=CELEX%3A32017R0745](https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX%3A32017R0745).
- 1034 59. Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April
1035 2016 on the protection of natural persons with regard to the processing of personal data
1036 and on the free movement of such data, and repealing Directive 95/46/EC (General
1037 Data Protection Regulation). [https://eur-lex.europa.eu/legal-](https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:32016R0679)
1038 [content/EN/TXT/?uri=CELEX:32016R0679](https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:32016R0679).
- 1039 60. Deng, J. *et al.* ImageNet: A large-scale hierarchical image database. in *2009 IEEE*
1040 *Conference on Computer Vision and Pattern Recognition* 248–255 (2009).
- 1041 61. Huang, G. B., Ramesh, M., Berg, T. & Learned-Miller, E. *Labeled Faces in the Wild: A*
1042 *Database for Studying Face Recognition in Unconstrained Environments*. [http://vis-](http://vis-www.cs.umass.edu/lfw/)
1043 [www.cs.umass.edu/lfw/](http://vis-www.cs.umass.edu/lfw/). (2007).
- 1044 62. Boyadjiev, S. A. & Jabs, E. W. Online Mendelian Inheritance in Man (OMIM) as a
1045 knowledgebase for human developmental disorders. *Clin. Genet.* **57**, 253–266 (2000).
- 1046 63. den Dunnen, J. T. *et al.* HGVS Recommendations for the Description of Sequence
1047 Variants: 2016 Update. *Hum. Mutat.* **37**, 564–569 (2016).

- 1048 64. Stevens-Kroef, M., Simons, A., Rack, K. & Hastings, R. J. Cytogenetic Nomenclature
1049 and Reporting. in *Cancer Cytogenetics: Methods and Protocols* (ed. Wan, T. S. K.) 303–
1050 309 (Springer New York, New York, NY, 2017).
- 1051 65. Kaye, J. *et al.* Dynamic consent: a patient interface for twenty-first century research
1052 networks. *Eur. J. Hum. Genet.* **23**, 141–146 (2015).
- 1053 66. Nellåker, C. *et al.* Enabling Global Clinical Collaborations on Identifiable Patient Data:
1054 The Minerva Initiative. *Front. Genet.* **10**, 611 (2019).
- 1055 67. Schoeman, L., Honey, E. M., Malherbe, H. & Coetzee, V. Parents' perspectives on the
1056 use of children's facial images for research and diagnosis: a survey. *J. Community*
1057 *Genet.* **13**, 641–654 (2022).
- 1058 68. Schmidt, A. *et al.* Next-generation phenotyping integrated in a national framework for
1059 patients with ultra-rare disorders improves genetic diagnostics and yields new molecular
1060 findings. *medRxiv* 2023.04.19.23288824 (2023) doi:10.1101/2023.04.19.23288824.
1061