

# **GestaltMML: Enhancing Rare Genetic Disease Diagnosis through Multimodal Machine Learning Combining Facial Images and Clinical Texts**

Da Wu<sup>1</sup>, Jingye Yang<sup>1</sup>, Cong Liu<sup>2</sup>, Tzung-Chien Hsieh<sup>3</sup>, Elaine Marchi<sup>4</sup>, Justin Blair<sup>5</sup>, Peter Krawitz<sup>3</sup>, Chunhua Weng<sup>2</sup>, Wendy Chung<sup>6</sup>, Gholson J. Lyon<sup>4,7</sup>, Ian D. Krantz<sup>5</sup>, Jennifer M. Kalish<sup>5,8,9</sup>, Kai Wang<sup>1,10\*</sup>

1 Raymond G. Perelman Center for Cellular and Molecular Therapeutics, Children's Hospital of Philadelphia, Philadelphia, PA 19104, USA

2 Department of Biomedical Informatics, Columbia University Irving Medical Center, New York, NY 10032, USA

3 Institute for Genomic Statistics and Bioinformatics, University Hospital Bonn, Rheinische Friedrich-Wilhelms-Universität Bonn, Bonn, Germany

4 Department of Human Genetics, New York State Institute for Basic Research in Developmental Disabilities, Staten Island, NY, USA

5 Division of Human Genetics, Children's Hospital of Philadelphia, Philadelphia, PA 19104, USA

6 Department of Pediatrics, Boston Children's Hospital, Harvard Medical School, Boston, MA, USA

7 Biology PhD Program, The Graduate Center, The City University of New York, New York, United States of America

8 Department of Genetics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, United States

9 Department of Pediatrics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, United States

10 Department of Pathology and Laboratory Medicine, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA 19104, USA

\*: correspondence should be addressed to wangk@chop.edu.

## **ABSTRACT**

Individuals with suspected rare genetic disorders often undergo multiple clinical evaluations, imaging studies, laboratory tests and genetic tests, to find a possible answer over a prolonged period of time. Addressing this “diagnostic odyssey” thus has substantial clinical, psychosocial, and economic benefits. Many rare genetic diseases have distinctive facial features, which can be used by artificial intelligence algorithms to facilitate clinical diagnosis, in prioritizing candidate diseases to be further examined by lab tests or genetic assays, or in helping the phenotype-driven reinterpretation of genome/exome sequencing data. Existing methods using frontal facial photos were built on conventional Convolutional Neural Networks (CNNs), rely exclusively on facial images, and cannot capture non-facial phenotypic traits and demographic information essential for guiding accurate diagnoses. Here we introduce GestaltMML, a multimodal machine learning (MML) approach solely based on the Transformer architecture. It integrates facial images, demographic information (age, sex, ethnicity), and clinical notes (optionally, a list of Human Phenotype Ontology terms) to improve prediction accuracy. Furthermore, we also evaluated GestaltMML on a diverse range of datasets, including 528 diseases from the GestaltMatcher Database, several in-house datasets of Beckwith-Wiedemann syndrome (BWS, over-growth syndrome with distinct facial features), Sotos syndrome (overgrowth syndrome with overlapping features with BWS), NAA10-related neurodevelopmental syndrome, Cornelia de Lange syndrome (multiple malformation syndrome), and KBG syndrome (multiple malformation syndrome). Our results suggest that GestaltMML effectively incorporates multiple modalities of data, greatly narrowing candidate genetic diagnoses of rare diseases and may facilitate the reinterpretation of genome/exome sequencing data.

### **Keywords:**

Multimodal Machine Learning, Artificial Intelligence, Large Language Models, Human Phenotype Ontology, Rare Genetic Disorders, Facial phenotyping

## INTRODUCTION

A substantial proportion of the global population, more than 6%, is affected by a rare genetic disorder<sup>1</sup>. While collectively common, rare diseases are individually rare<sup>2</sup>. Rare diseases are typically defined as affecting fewer than 200,000 people in the USA or less than one in 2,000 of the general population in Europe<sup>3</sup>. Based on the latest Orphanet<sup>4</sup> and OMIM<sup>5</sup> databases, currently there are at least 7000 rare genetic diseases. Due to the inherent rarity and extensive phenotypic heterogeneity of rare genetic disorders, accurately making a genetic diagnosis presents a challenge, often leading to a “diagnostic odyssey”<sup>6-8</sup>. Patients with suspected genetic syndromes often need to undergo multiple clinical evaluations, imaging studies, and laboratory tests, in addition to multiple modalities of genetic tests, including karyotype, chromosome microarray, gene panels, exome sequencing or genome sequencing, to make the diagnosis. Clinicians often encounter difficulties deciding what diagnostic test to use for efficient and accurate diagnosis, as they must navigate long differential diagnoses for each of many different symptoms. Shortening the odyssey could have significant clinical, psychosocial, and economic benefits<sup>8,9</sup>.

Many genetic diseases have distinctive facial features or dysmorphism (collectively considered the “facial gestalt”), which often provide important clues on facilitate diagnoses and expedite referrals to domain experts or suggest targeted genetic tests. In some cases, the recognition of a syndrome from a facial gestalt can be the first step in making a diagnosis<sup>10</sup>. However, the effectiveness of facial recognition relies heavily upon the clinician's experience with facial recognition of syndromes. Given the many hundreds of rare genetic diseases with facial dysmorphisms, some only identified in the last 5 years, the facial recognition task is prohibitive for any clinician.

Following the recent success in Computer Vision (CV), there are several next generation phenotyping (NGP) approaches developed to analyze and predict rare genetic disorders based on patient's 2D frontal facial images<sup>11-13</sup>. Among those, one widely known approach is called DeepGestalt<sup>11</sup>, which was developed by FDNA Inc. as the “Face2Gene” product, and was pretrained on deep convolutional neural network (DCNN) using CASIA<sup>14</sup> and later fine-tuned on over 17106 patient frontal facial images with 216 disorders. However, DeepGestalt was only trained on a limited number of syndromes, which accounts for only a small proportion of the total syndromes. Adding a newly discovered syndrome requires collecting and adding new images and retraining the model. To make the model more inclusive for new “unseen” syndromes, GestaltMatcher<sup>12</sup> was introduced as an improvement to DeepGestalt, in the sense that it takes the DeepGestalt's feature layer before the final classification layer as a common embedding space (also known as “Clinical Face Phenotype Space” - CFPS) to encode learned facial dysmorphic features. Every frontal facial image was encoded into a 320-dimensional feature representation vector and by doing this, it can quantify the distance between different images and further identify the “closest match” among patients with known or unknown disorders, regardless of prevalence. Furthermore, one additional advantage of GestaltMatcher is that there is no need to alter the model's architecture and retrain the model, when integrating newly identified syndromes. Despite these aforementioned successes, both DeepGestalt and GestaltMatcher use relatively dated model architecture and datasets for transfer learning introduced by Yi et al.<sup>14</sup> More recently, Hustinx *et al.*<sup>13</sup> updated the model architecture with more advanced iResNet<sup>15</sup> and ArcFace<sup>16</sup> and used various updated facial image datasets, including VGG2<sup>17</sup>, CASIA<sup>14</sup>, MS1MV2<sup>16</sup>, MS1MV3<sup>16</sup> and Glint360K<sup>18</sup>, for pretraining. They also tried on different loss functions and proposed a model ensemble (combining three ArcFace

models) to integrate face verification and disorder-specific models to improve performance on both seen and unseen syndromes. The model ensemble can achieve higher accuracy on unseen syndromes than all the previous models after fine-tuning.

Nonetheless, in numerous instances, facial images alone do not provide adequate information to make a precise diagnosis. For instance, syndromes such as Noonan syndrome (NS), Prader-Willi syndrome (PWS), Silver-Russell syndrome (SRS), and Aarskog-Scott syndrome (ASS) all have severe to moderate short stature<sup>19</sup> which cannot be effectively reflected in frontal facial pictures. Additional phenotypic traits such as sleep disturbances, impaired balance and intellectual disability cannot be effectively captured by facial or other body photos. These aspects require additional data types (e.g., clinical notes). Moreover, numerous investigations<sup>20-28</sup> have been conducted examining the contribution of age, sex, as well as racial and ethnic differences, to the phenotypic expression and frequency of various disorders and syndromes. Certain groups, often categorized as minorities, encounter challenges that stem from systemic biases ingrained within data availability, collection and analysis processes. These biases can inadvertently lead to misrepresentations, inaccuracies, and disparities in the rare genetic disorder predictions of these groups. Motivated by those facts, there are already some models developed trying to integrate facial images and clinical HPO terms together. The authors introduced the “prioritization of exome data by image analysis” (PEDIA) strategy<sup>29</sup>. PEDIA incorporates sequence variant interpretation with insights from the advanced phenotyping tool DeepGestalt. This approach enhances clinical assessments by combining expert human evaluation and artificial intelligence analysis, using frontal photographs to provide a more comprehensive assessment of individual clinical presentations. More recently, the PhenoScore<sup>30</sup>, an AI-based framework for analyzing genetic syndromes, was introduced and it comprises two modules: facial feature extraction from 2D photographs and HPO-based phenotypic similarity calculation. The framework uses a trained Support Vector Machine (SVM) for syndrome classification, based on extracted facial features and HPO similarities. However, these existing models process images and texts separately, then combine the results. This type of approach to integrating multi-modality data may lead to information loss, as it fails to fully capture the interaction between different modalities during training and it uses *ad hoc* methods to assign weight and combine information. In addition to the advancements mentioned earlier, a recent development in the field is DxGPT<sup>31</sup>, a text-only GPT-based model tailored for diagnosing rare genetic diseases. This model is built upon the closed-source GPT-4. In light of this, our objective is to create a multimodal machine learning (MML) methodology that incorporates a sophisticated modality interaction module. This methodology will handle both facial images and clinical texts in a uniform manner. The intended methodology aims to effectively merge patient facial images with textual information, which includes demographic details such as age, gender, and ethnicity, along with clinical notes, thereby preserving the integrity and richness of the data.

The recent progress in Transformer-based multimodal machine learning models has made our objective attainable. The story of Transformers started with the landmark paper “Attention is all you need”<sup>32</sup> which introduces the so-called *self-attention mechanisms*, enabling the model to process sequence (e.g. texts sentences) in parallel rather than sequentially, like the traditional recurrent and convolutional neural networks. This design is revolutionary in the sense that it leads to improved performance, faster training, and scalability (larger size leads to better performance). Since then, Transformers have been extensively applied to both Natural Language Processing (NLP) and Computer Vision (CV), showcasing their versatility and

effectiveness in various tasks. For instance, in NLP, Transformers have been applied to machine translation<sup>33</sup>, text generation<sup>34</sup>, sentiment analysis<sup>35</sup>, named entity recognition (NER)<sup>36</sup>, and others. In CV, tasks such as image classification<sup>37</sup>, object detection<sup>38</sup>, image segmentation<sup>39</sup>, image captioning<sup>40</sup>, visual question answering (VQA)<sup>41</sup> now all rely on the Transformers. Recently, several multimodal machine learning models that leverage the strengths of Transformers have been developed, for instance, ViLT (Vision-and-Language Transformer Without Convolution or Region Supervision)<sup>41</sup>, CLIP (Contrastive Language–Image Pre-training)<sup>42</sup>, VisualBERT<sup>43</sup>, ALBEF (Align Before Fuse: Vision and Language Representation Learning with Momentum Distillation)<sup>44</sup> and Google Gemini<sup>45</sup>.

Taking all these factors and tools into consideration, for the task of predicting rare genetic disorders, we introduce a novel methodology, *GestaltMML*, utilizing the ViLT (Vision-and-Language Transformer)<sup>41</sup>, which has the *simplest* design among vision-and-language models as it uses the nontrivial Transformer module for modality interaction learning while only using trivial (linear) vision and textual embedding.

## RESULTS

### **Summary of the computational experiments**

The overall workflow of the study is illustrated in **Fig 1A**. For *GestaltMML*, the data pre-processing procedure is illustrated in **Fig 1B** with an example encompassing both text and image data. The transformer model architecture is summarized in **Fig 1C**. In the following sections, we first demonstrate the performance of *GestaltMML* on the GestaltMatcher Database (GMDB)<sup>12,46</sup>. *GestaltMML* is a multimodal machine learning model that integrates facial images with demographic information and clinical HPO texts. GMDB is a collection of curated medical photographs of genetic syndromes as a resource for clinician and computer scientists. As shown in **Table 1**, the version of database (v1.0.9) that we used in this study contains 9764 frontal facial images from 7349 patients affected with 528 rare genetic disorders.

The database includes patients of diverse ancestry through global collaboration. The specific ancestral categories represented are Middle-East/West Asian, American – Native, South-East Asian, North African, Unknown, African American, American - Latin/Hispanic, East Asian, Asian Others, South Asian, Others, Sub-Saharan, and African. **Fig. 2A** shows a significantly skewed distribution, with 59.48% of patients being of European ancestry. As highlighted in the recent efforts on constructing the diverse GMDB database<sup>46</sup>, despite significant efforts to create as diverse a dataset as possible, achieving complete balance in patient characteristics is challenging due to the nature of rare diseases. This imbalance introduces inevitable difficulties for AI models in diagnosing rare genetic diseases, hence the decision to include ethnicity information in the textual data. **Fig. S1** illustrates a fairly equal distribution between males and females. There is an uneven age distribution, with a majority (64.90%) of patients being under 5 years old.

Following the convention in previous work on developing the ensembled image model<sup>13</sup>, for some evaluations, we measured performance on the GMDB-frequent and GMDB-rare subset, based on the number of patients for each disease; The GMDB-frequent contains disorders with more than ( $>$ ) 6 patients, while GMDB-rare contains disorders with less or equal to ( $\leq$ ) 6 patients. Additionally, we will explore the significance of text and image features in *GestaltMML* and compare its effectiveness against current image-based models.

Finally, we extended our evaluation to several external validation datasets, encompassing patient data from the Children’s Hospital of Philadelphia (CHOP), New York State Institute for Basic Research in Developmental Disabilities (NYSIBRDD), and published literature. Our model also exhibits high performance on some of these external datasets, demonstrating the robustness of our methods.

### ***GestaltMML accurately classifies rare genetic diseases in GMDB***

To thoroughly assess the predictive capabilities of the image-text pairs and overcome the challenge posed by a substantial amount of missing text data in cross-validation experiments, we constructed a train-test split, called *optimal train-test split*, in which the test set comprises images that meet two criteria: (1) they possess non-null present features and (2) they exclusively represent disorders that the training dataset has encountered with non-null features. For detailed algorithms for constructions of optimal train-test split with x:1 train: test ratio, see subsection “*Construction of optimal train-test splits*” in section “*Methods*.”

In the end, we constructed optimal train-test splits for each optimal train: test ratio from 1:1 to 9:1. The details on the number of images were summarized in **Table S1**. To make our results more robust, we constructed the optimal train-test splits using three different random seeds and calculated the means and standard deviations for all Top-1, Top-10, Top-50, Top-100 accuracies. Top-N accuracy is defined as the measure of precision when the actual disease label is included within the top-N predictions, whether from image information alone, or from the combined image and textual information. With optimal train: test ratio of 3:1, the model can reach mean accuracies of 72.54% for Top-1, 83.59% for Top-10, 88.96% for Top-50, and 91.64% for Top-100. The mean accuracies and their corresponding standard errors for other optimal train: test ratios are illustrated in **Table S2**.

By incorporating all the available images, demographic information, and clinical features, we observe that even with 3:1 train-test ratio, the model can reach high testing accuracy. While the accuracy is impressive, we acknowledge that the number of diseases that are analyzed is relatively small (528 in GMDB compared to the thousands of known rare diseases). Furthermore, there are likely ascertainment biases in typical facial image databases such as GMDB, such that only diseases for which there are characteristic dysmorphic features are documented, so the method may not work well for rare diseases without machine-recognizable facial features. Nevertheless, the use of demographic information and clinical phenotype information in GestaltMML can facilitate the prioritization of diseases in these cases.

### ***Feature importance analysis and comparisons with existing image models on GMDB dataset***

Until now, almost all the existing literature has primarily centered on employing facial images alone for the prediction of rare genetic disorders. We selected the up-to-date ensembled image model documented in the previous work<sup>13</sup> as our initial benchmark model. We subsequently conducted an extensive comparative analysis. This included our modified versions of GestaltMML, aimed at exploring feature importance, discussed below. The summary of these comparisons can be found in **Table 2**. It is worth noting that GestaltMML *exclusively* employs the Transformer architecture and are completely devoid of convolution. In stark contrast, none of the other image-only models listed in previous work of<sup>13</sup> utilize the Transformer architecture.

To make fair comparisons, we adopted the same train-test splits as in the previous work of developing ensembled image model<sup>13</sup>. Following the conventions in there<sup>13</sup>, we first separate the GMDB into GMDB-frequent and GMDB-rare. Recall that the GMDB-frequent contains disorders with more than 6 patients, while GMDB-rare contains disorders with less or equal to 6 patients. As described in **Table 1** there are 8547 images from 6376 patients with 244 disorders in GMDB-frequent and 1217 images from 973 patients with 284 disorders in GMDB-rare. For GMDB-frequent, we used 7755 images for training and the remaining 792 images for testing. For GMDB-rare, we used the same 10-fold cross validation as in the previous work of developing ensembled image model<sup>13</sup>. On average, there are 856.9 images for training and the remaining 360.1 for testing. We fine-tuned the training data from combined GMDB-frequent and GMDB-rare and further evaluate on GMDB-frequent or GMDB-rare. Altogether, the training involves 8611.9 images, with 792 images designated for testing in GMDB-frequent and 360.1 images for testing in GMDB rare.

To delve deeper into the significance of each modality (images and texts) and understand their individual effects on the final prediction power, we conducted an evaluation procedure that we call “modality-masking.” To test the prediction power of images, we masked out all the text. We fine-tuned the training data with entire texts components replaced by “\*” (facial image + \*) on ViLT. We call the fine-tuned model GestaltViT to emphasize that during the fine-tuning process, we exclusively focus on fine-tuning the facial images and not the image-text pairs. Similarly, we assessed the predictive capabilities of text by excluding all images. Specifically, we fixed one patient photo for all training samples, while keeping all other aspects unchanged as described earlier. The results obtained from this test were solely evaluated on images containing existing HPO terms. The fine-tuned model is denoted as GestaltLT, with “LT” standing for “Language Transformer”. Noticeably, as shown in **Table 2**, GestaltViT exhibits poorer performance compared to the ensembled image model. This outcome is entirely expected due to the ViLT model's linear patch embedding on the image parts and it aligns with the well-known *scaling property* of Transformer-based models, as discussed and compared in prior works like<sup>37</sup>. It is reasonable to anticipate that with a larger training dataset, the performance of GestaltViT would further improve. Regarding GestaltLT, it exhibits only slightly inferior performance compared to GestaltMML and outperforms both GestaltViT and the ensembled image model. Despite the natural biases present in the textual data due to GMDB and the use of data augmentation via OMIM, it exhibits remarkable predictive capabilities.

### ***GestaltMML demonstrates enhanced diagnostic equity for patients from underrepresented groups***

GestaltMML underwent training using the GMDB (v1.0.9), which encompasses data from several ethnically under-represented groups, identified as “Middle-East/West Asian,” “American – Native,” “South-East Asian,” “North African,” “Unknown,” “African American,” “American - Latin/Hispanic,” “East Asian,” “Asian Others,” “South Asian,” “Others,” “Sub-Saharan,” and “African.” In **Fig. 2B**, the mean accuracies of GestaltMML are showcased across various modalities used for inference. It is evident that clinical texts exert the most significant influence in enhancing performance, while demographic information also proves beneficial in augmenting results, particularly for minority patients. **Fig. 2C** illustrates that GestaltMML, by integrating frontal facial images, demographic details, and clinical texts, significantly enhances its predictive accuracy across under-represented ethnic groups, particularly when compared to training

exclusively on individuals of European descent - with only rare exceptions. Further elaboration on the methodologies for dividing the training and testing sets, as well as the detailed training and testing processes, is provided in the Methods section. Complete results, encompassing Top-1, Top-10, Top-50, and Top-100 mean percentage accuracies for each individual ethnicity, along with their corresponding standard deviations, are documented in **Tables S4-S9**.

### ***GestaltMML performs well on external validation datasets***

#### Overview

Although GestaltMML has demonstrated success on the GMDB database, we aim to evaluate its performance further on external validation data to assess its resilience to potential bias inherent in the GMDB database. Several case studies, including Beckwith-Wiedemann syndrome (BWS), Sotos syndrome, NAA10 neurodevelopmental syndrome, Cornelia de Lange syndrome (CdLS) and KBG syndrome are described and discussed separately below.

In this section, we use GestaltMML trained on optimal training: test split with ratio 9:1 (**Table S1**). However, for NAA10 patients, we employed GestaltMML that was trained using GMDB v1.0.3, with an optimal train: test ratio of 4:1, to avoid train-test overlap (the NAA10 patients from external validation cohort were included in the recent database update). Furthermore, we also compare the performance of GestaltMML with the state-of-the-art ensembled image model developed in the previous work<sup>13</sup>.

#### Beckwith-Wiedemann Syndrome (BWS)

Beckwith-Wiedemann Syndrome (BWS) is one of the most common overgrowth syndromes<sup>47-50</sup>. It exhibits both genetic/epigenetic and clinical diversity<sup>51,52</sup>. BWS involves molecular aberrations within a cluster of imprinted genes on chromosome 11p15.5-11p15.4: Loss of methylation at imprinting control region 2 (IC2) on the maternal allele is found in about 50% of patients, while paternal uniparental isodisomy for the 11p15 region (pUPD11) occurs in about 20% of patients<sup>53</sup>. Key signs of overgrowth like macrosomia, organomegaly, and hemihypertrophy, which may not be adequately represented in frontal facial images<sup>52,54,55</sup>. Consequently, earlier image-based models may fall short in providing reliable diagnosis when these important clinical features are not considered. To address this limitation, a multimodal machine learning approach that incorporates clinical texts (in the form of HPO terms) is needed.

We collected two groups of in-house patients affected with BWS at CHOP: one group with imprinting control region 2 loss of methylation (IC2) and another with paternal uniparental isodisomy of chromosome 11 (pUPD11). Each patient's clinical phenotype is described, and all but one had complete demographic details, including sex, ethnicity, and age. We applied our GestaltMML to these data, with the findings presented in **Table 3**. The results show that clinical phenotypic descriptions are highly valuable, as GestaltMML achieves 100% detection Top-1 accuracy with this information. Nonetheless, the accuracy of GestaltMML markedly diminishes when it depends exclusively on facial images. Similarly, even the state-of-the-art ensembled image model struggles to accurately detect BWS. To illustrate, we showcase a particular example in **Fig. 3A**, where facial image information alone ranks BWS as the sixth most likely diagnosis yet adding demographic information and clinical phenotype descriptions improves the rank of the disease to the first.



### Sotos Syndrome

Sotos syndrome is another rare genetic disorder characterized by excess growth during the early years of life, and children with Sotos syndrome typically have greater height, weight, and larger head size (macrocephaly) compared to their peers<sup>56-59</sup>. Sotos syndrome frequently involves delays in motor skills, cognitive abilities, and social development. Since many of the clinical phenotypic features cannot be represented by facial photos, this syndrome is tested here to examine the importance of employing clinical texts for effective and accurate diagnosis, and to also compare with another overgrowth syndrome BWS.

We gathered data from 23 patients with Sotos Syndrome at the CHOP and conducted predictions using GestaltMML. The findings are detailed in **Table 3**, where again it is evident that incorporating multimodal data, including demographic details and clinical texts, greatly enhances the accuracy of inference. A particular instance of this is demonstrated in **Fig. 3B**, where facial image information alone ranks Sotos syndrome as the 288<sup>th</sup> likely diagnosis yet adding demographic information and clinical phenotype descriptions improves the rank of the disease to the first. We also discovered that GestaltMML most frequently misdiagnoses Sotos syndrome as Marshall-Smith Syndrome (OMIM:602535), a genetic disorder characterized by distinctive facial traits such as prominent forehead, shallow orbits, blue sclerae, depressed nasal bridge, and micrognathia<sup>60</sup>.

### NAA10-related Neurodevelopmental Syndrome

NAA10-related neurodevelopmental syndrome<sup>61</sup> is an X-linked condition with a broad spectrum of findings ranging from a severe and often lethal phenotype cardiac in males (five deceased boys)<sup>62</sup>, to the severe NAA10-related intellectual disability in both males and females. In 2023, we expanded the phenotypic spectrum of NAA10-related neurodevelopmental syndromes through analysis of 56 individuals with NAA10 variants, demonstrating a phenotypic spectrum that includes variable intellectual disability, delayed milestones, autism spectrum disorder, craniofacial dysmorphism, cardiac anomalies, seizures, and visual abnormalities<sup>63</sup>. We collected clinical information (photos and clinical texts) on 68 subjects from NYSIBRDD. Note that they are not included in the previous version of the GMDB, i.e., v1.0.3, but most of them are now included in the new GMDB (v1.0.9). Therefore, results in this section are based on trained model on GMDB (v1.0.3) only.

We used the same testing procedure as used earlier using GestaltMML. Regarding the text data, we extracted demographic information of patients and HPO terms from the clinical summaries provided by NYSIBRDD. **Fig. 4** illustrates the ranking of true label among a total of 449 disease labels (total number of labels in GMDB v1.0.3), comparing results obtained from facial images alone with those derived from a combination of facial images, demographic information, and clinical phenotype descriptions. In almost all the cases, the use of multimodal information improves the prediction accuracy for GestaltMML significantly. In some instances, incorporating textual information can lead to poorer prediction outcomes. This primarily stems from the similarity of the text component in our test data to those of other neurodevelopmental disease labels in GMDB. Disease labels such as Intellectual Developmental Disorder, X-Linked, Syndromic 33 (OMIM:300966) and Developmental Delay, Hypotonia, Musculoskeletal Defects,

and Behavioral Abnormalities (OMIM: 619595) have comparable textual descriptions. Such similarities can cause great confusion in the model and consequently degrade the results.

Additionally, we stress that demographic information can introduce certain biases in prediction outcomes. This bias is largely due to the uneven representation of demographic groups in the training set. For nearly all the diseases, most patients are under 5 years old and of European ancestry. When demographic information of a new patient falls outside the typical range in GMDB, the discrepancy is likely to result in unstable predictions. In future research, we intend to integrate more diverse data sources to improve the predictive power of GestaltMML, particularly for neurodevelopmental syndromes.

### Cornelia de Lange syndrome (CdLS)

Cornelia de Lange syndrome (CdLS), also known as Brachmann-de Lange syndrome, is a genetically heterogeneous multiple malformation syndrome typically characterized by growth restriction, variable upper limb differences, hypertrichosis, long eyelashes, thick eyebrows, short nasal root and tip with anteverted nares, long philtrum and thin upper lip and other findings<sup>64-66</sup>.

The external validation dataset for CdLS patients, collected from CHOP with 19 samples, underwent the same evaluation through both GestaltMML and an ensemble image model, as detailed in **Table 3**. Contrary to previous conditions, CdLS diagnosis favored image-based over text-based analyses. We found that GestaltMML, even when utilizing only facial images, can surpass the performance of multimodal inference. Therefore, it is not unexpected that the ensembled image model can attain exceptionally high prediction accuracy. This is partially attributed to CdLS's distinctive facial features that significantly aid in accurate diagnosis. However, textual data, while informative with details such as “global developmental delay” and “feeding difficulties,” tend to blur distinctions with other neurodevelopmental syndromes in GMDB (v1.0.9), impacting the model's accuracy. Our analysis of this syndrome with well known facial features illustrate that GestaltMML may be more useful for other syndromes with subtle features that are hard to recognize by human experts.

### KBG syndrome

KBG syndrome is an extremely rare, pan-ethnic, autosomal dominant disorder characterized by macrodontia, post-natal short stature, skeletal anomalies, abnormal hair implantation, and developmental delays<sup>67-71</sup>.

The external validation cohort for KBG syndrome from NYSIBRDD, comprising 18 samples, was evaluated using the same testing methodology. Results in **Table 3** again indicated superior performance of image-based models over both multimodal and single-text models for KBG syndrome. Likewise, as with the previous case, this disparity is partly attributed to the older age of patients in the outside validation set compared to the training set, where encoding age as text introduced prediction instability. Additionally, same as the case of CdLS, the presence of common HPO terms in clinical texts, like “global developmental delay” and “intellectual disability, severe”, which overlap with many other rare diseases in GMDB (v1.0.9), contributed to predictions towards other neurodevelopmental syndromes despite the additional information provided by text data. This case study again highlighted the importance of having a training facial photo database with a large range of age distributions.

### ***GestaltMML exhibits outstanding performance in clustering diseases that have clinical similarities.***

Finally, we performed a two-component UMAP clustering analysis on the logit values from the penultimate layer of the GestaltMML model (see Methods). This analysis focused on three comparative sets of diseases: BWS versus Sotos Syndrome, NAA10 versus NAA15-related syndromes, and KBG Syndrome versus Cornelia de Lange Syndrome (CdLS).

The first set of analysis focuses on external validation data including two BWS patient cohorts (IC2 and pUPD11, as previously discussed) in conjunction with patients with Sotos syndrome from previous section. While both diseases represent overgrowth syndromes, the model can clearly separate these two, and even separate the two genetic subtypes of the same syndrome (**Fig. 5A**).

Next, within the GMDB (v1.0.9), we evaluate GestaltMML's clustering efficiency for patients associated with NAA10 and NAA15-related neurodevelopmental syndromes (**Fig. 5B**). This analysis was confined to patients cataloged in GMDB (v1.0.9) only. We found that despite overlap of clinical phenotypes between these two syndromes, the model can still separate those affected with NAA10 deficiency versus NAA15 deficiency.

Lastly, we conducted tests using external validation data for patients with KBG syndrome and Cornelia de Lange syndrome (CdLS), as mentioned in the previous section (**Fig. 5C**). As expected, these two diseases can be separated. However, it is worth noting that among CdLS patients, facial image inference reveals two distinct clusters. Additional investigation found that this occurrence is attributed to background color variations: one cluster comprises images with white or pale backgrounds, while the other consists of images with warm or dark yellow backgrounds. This observation suggests that additional improvements to normalize background color may increase precision of the representation of facial images.

## **DISCUSSION**

In the current study, we introduced a novel multimodal approach, GestaltMML, to integrate frontal facial photos, clinical features, and demographic information together to narrow the differential diagnosis of rare genetic diseases. This approach is motivated by the observation that sole reliance on facial images of patients is insufficient to encompass all essential information required for the accurate diagnosis of rare genetic disorders. Our findings indicate that multimodal machine learning can lead to a substantial enhancement in the accuracy of predicting likely genetic diagnoses. Furthermore, it proves to be an indispensable resource for differentiating rare disorders with shared clinical characteristics via UMAP clustering analysis. Similar to the capabilities of the previously developed GestaltMatcher<sup>12</sup> and the Ensembled image model<sup>13</sup>, this approach to clustering enables the model to automatically identify novel, previously unrecognized rare diseases without the need to alter the classification layers or undergo a complete retraining of the model. In combination with genome/exome sequencing data, GestaltMML is likely to greatly facilitate the interpretation or periodic reinterpretation of data, ultimately addressing the “diagnostic odyssey” challenge.

As previously mentioned, GestaltMML utilizes only the Transformer architecture. This choice aligns with the foundational principles outlined in the seminal paper “Attention is all you need,”<sup>32</sup>

which advocates for the complete substitution of recurrent or convolutional networks. Compared to classical CNN-based image models, additional technical differences and key innovations of our approach are discussed as follows: (1) Our methodology diverges from prior models that focused solely on facial images by incorporating both facial images and texts as inputs for the prediction of rare genetic disorders. This distinction sets it apart from the image-only models discussed in the previous work<sup>13</sup> and the related references therein. (2) We integrated demographic data from patients, including sex, age, and race/ethnicity details, into the text inputs, enabling the model to discern distinct patterns for each rare disorder. We demonstrated that this approach successfully reduces biases inherent in data collection and analysis, especially regarding underrepresented minority groups, leading to a fairer diagnostic procedure. (3) We introduced a data augmentation technique leveraging the OMIM database<sup>72</sup>. This approach enhances the model's training process by infusing it with a rich and comprehensive textual knowledge base. The incorporation of information from the OMIM database during training contributes to improved performance of the model. (4) We further examined the significance of textual and visual elements during multimodal training using *modality masking* techniques, offering valuable insights for future research endeavors. (5) Ultimately, upon evaluating GestaltMML and the ensembled image model across several external validation datasets, we observed a notable improvement in diagnostic accuracy for numerous conditions, such as BWS and Sotos syndrome, through the integration of textual information. Conversely, for diseases like CdLS and KBG syndrome, image-only models (this includes both the image segment of GestaltMML and the ensemble of image models) outperformed multimodal methods in terms of prediction effectiveness. These findings prompt clinicians to be judicious in using multimodal GestaltMML to assist in diagnosis, particularly when clinical HPO terms resemble those of many different disorders. The feature importance analysis of GestaltMML indicates that the text component's predictive power exceeds that of images, suggesting that non-specific HPO texts may confuse the model and decrease its accuracy. Where facial images are distinctly recognizable, basing a diagnosis solely on these images might yield higher precision.

To optimize performance, the GestaltMML model leverages a straightforward approach of using concatenated HPO terms as textual inputs. Incorporating continuous clinical text paragraphs directly into these models may, however, affect the performance of GestaltMML adversely. The latest developments in large language models (LLMs) have significantly enhanced the capability to identify and extract HPO terms directly from clinical text paragraphs with high efficiency. A notable example is PhenoGPT<sup>73</sup>, which, built upon advanced large language models, demonstrates high accuracy in extracting HPO terms. Given its effectiveness, it is advisable to preprocess continuous clinical text paragraphs with PhenoGPT or similar large language models before feeding it into GestaltMML for optimal results.

Notwithstanding the several advantages and strengths outlined so far, it is important to acknowledge that our current GestaltMML methodology does have limitations. Here we highlight those limitations: (1) The major limitation pertains to image embedding within our GestaltMML framework. Our current approach employs linear patch embedding (same as ViT<sup>37</sup>), a method that has demonstrated comparably lower efficacy when compared to textual embedding, particularly when working with constrained training data. This limitation of the image modules also prompts caution regarding multimodal inference, particularly when text components lack distinctiveness, such as the cases of CdLS and KBG syndrome discussed before. To address this concern, we recommend the adoption of a more sophisticated feature extraction module, such as those based on Transformers or Convolutional Neural Networks (CNNs). For more

comprehensive discussion, see literatures<sup>42,43</sup>. (2) The GestaltMML is constructed upon the foundations of ViLT, a pretraining framework designed for both images and text, albeit not specific to facial images with medical contexts. There is an intriguing avenue to explore: the creation of a foundational multimodal model that is pretrained using facial images (ideally from patients) alongside corresponding clinical textual descriptions. However, we are aware of the challenges to procure a large and diverse dataset encompassing patient facial images and medical captions, especially since much of the data are not in the public domain or are not consented for research use. Moreover, the training process for such a model might demand a substantial investment of time and computing resources. With the expansion of facial photo databases such as GMDB as well as the integration of photo information to clinical phenotype databases on rare diseases, it may be possible to create a foundational multimodal model down the road.

## **METHODS**

### ***Patients and Photos***

The study to develop multimodal machine-learning approaches for rare disease diagnosis was approved by the Institutional Review Board of the Children’s Hospital of Philadelphia (IRB 18-015712). The GMDB (v1.0.9) database used in the current study contains 9764 images from 7349 patients affected with 528 genetic disorders, obtained from <https://db.gestaltmatcher.org/>. BWS and Sotos syndrome images were collected and analyzed under the oversight of the Children’s Hospital of Philadelphia (CHOP) Institutional Review Board protocol (IRB 13-010658). The collection and analysis of data on individuals with Cornelia de Lange syndrome were performed under CHOP Institutional Review Board protocol (IRB 16-013231). In brief, consent was obtained from all patients and/or legal guardians to analyze and in some case publish the images. For collection and facial phenotyping analysis on the NAA10-related neurodevelopment syndrome and KBG syndrome, both oral and written patient consent were obtained for research and publication, with approval of protocol #7659 for the Jervis Clinic by the New York State Psychiatric Institute - Columbia University Department of Psychiatry Institutional Review Board. Written family consent was given for publication of any photographs of the children.

### ***Training and Evaluation of GestaltMML***

#### *Overview of training data sources*

In the developments of GestaltMML, we will primarily use two sources of data. The first one is the GMDB database, which contains frontal facial images of patients and corresponding textual metadata. It is open to researchers in medical domains, and one needs to apply for access first to use the data. The second one is the OMIM website<sup>72</sup> which serves as our ground knowledge base for rare genetic disorders. To deal with large amounts of missing textual data in GMDB, we will use the textual data from OMIM database for the purpose of data augmentation. The visual representation of whole data preprocessing procedure can be found in **Fig. 1B**.

#### *Image data preprocessing*

Our training and test images were cropped by the open source “FaceCropper” described in the GitHub page<sup>13</sup>. In our experiment, the facial images are of dimension 112 \* 112. Alternatively,

the image can be manually resized to these dimensions, ensuring that the primary face encompasses the entire picture. Notice that the original facial images are subjected solely to cropping, without any alterations such as flipping, rotating, or converting to grayscale.

### Text data preprocessing

To make the textual data preprocessing procedure clearer, we will separate our discussions into two cases.

The first case concerns images that have non-null present features, or equivalently, at least one HPO id in the “present features” column of the metadata. In this case, we will do the following two steps: (1) Transform the HPO id(s) into real text data via the standard HPO dictionary<sup>74</sup> and then concatenate them with empty space in between. For instance, the “HP:0000486; HP:0001263; HP:0010864” will become “Strabismus Global developmental delay Intellectual disability severe.” (2) Add patients' demographic information in the front. The image metadata of GMDDB database contains patients' sex, age, and ethnicity (or ethnicity note), which will be combined for our model training. For instance, the demographic textual data will look like “Sex male Age 4 years 8 months Ethnicity European.” If there is missing information, then we will simply leave that space empty. Therefore, the combined textual data in our first case will look like “Sex male Age 4 years 8 months Ethnicity European Strabismus Global developmental delay Intellectual disability severe.”

The second case deals with the case when images do not have present features at all, or equivalently, no HPO id in the “present features” column of the metadata. In this case, we will do the following two steps: (1) Use the textual data in the “clinical features” section of OMIM database as the primary source for data augmentation. Due to the limitations of model inputs' length and for the sake of saving computing power and budgets, we further use OpenAI's ChatGPT<sup>75</sup> to summarize those texts within 500 tokens. The prompt we gave is “Summarize most crucial phenotype characteristics of the following texts describing clinical features of some rare genetic disorder within 500 tokens”. The sample texts paragraph (after summarization by ChatGPT) looks like “Clinical features of this rare genetic disorder include supravalvular aortic stenosis SVAS mental retardation distinctive facial features dental anomalies peripheral pulmonary artery stenosis infantile hypercalcemia statural deficiency characteristic dental malformation and a hoarse voice Other features may include renal abnormalities cardiovascular disease joint limitations hypotonia delayed growth cataracts stroke and cognitive deficits Patients often have musical and verbal abilities but struggle with visual-motor integration and attention deficit disorder They may also exhibit hypersensitivity to sounds and have urinary abnormalities.” Note that we also remove all the “,” “.” “:”, “()”, etc. to save token space for training. Same thing applies to the first case. (2) Likewise, add patients' demographic information in the front.

The text data pre-processing approach mentioned is tailored exclusively for the GMDDB database. Like mentioned before, for clinical practitioners working in real-world settings, it is recommended to employ PhenoGPT for the extraction of HPO terms from clinical text paragraphs. This should be followed by concatenating these terms with demographic data, as previously described, to ensure optimal data preparation and analysis.

### Construction of optimal train-test splits

The optimal  $x:1$  split in **Table S1** is constructed as follows: (1) Select all the disorders that have images with non-null present features and denote the set of such disorders  $D$ . (2) For each disorder  $d$  in  $D$ , let  $I_d$  denote the set of all the image ids of disorder  $d$  with non-null present features and compute the cardinality  $|I_d|$ . Next, randomly select  $\lfloor |I_d|/(x+1) \rfloor$  image ids from  $I_d$  and call them  $I_d^t$ . For instance, if  $|I_d| = 5$  and train:test ratio is 3:1, then  $\lfloor 5/(3+1) \rfloor = 1$  image will go to the test set  $I_d^t$  and rest 4 images will be grouped into training set. On the other hand, if  $|I_d| = 2$  and train:test ratio is 3:1, then  $\lfloor 2/(3+1) \rfloor = 0$  image will go to test set. In other words, we do not select any test image under this optimal train-test split ratio. (3) The total testing set is simply the union  $\cup_{d \in D} I_d^t$  and the training set is the complement of the testing set.

Using the above algorithm, we constructed optimal train-test splits for each optimal train: test ratio from 1:1 to 9:1. To make our results more robust, we will repeat the above algorithms three times (using three different random seeds) and calculate the means (**Table S1**) and standard deviations (**Table S2**) for all Top-1, Top-10, Top-50 and Top-100 accuracies.

#### Train-test splits for experiments of GestaltMML enhancing diagnostic equity for patients from minority groups.

In **Fig. 2B**, we assess the effect of only including demographic information and including both demographic information and clinical HPO terms in the textual component. The text component will look like (1) "\*" (no texts at all), (2) "Sex male Age 4 years 8 months Ethnicity European" (demographic information only) or (3) "Sex male Age 4 years 8 months Ethnicity European Strabismus Global developmental delay Intellectual disability severe" (demographic information and clinical HPO texts).

In **Fig. 2C**, we evaluate the diversity of training datasets that include both facial images and textual information. We begin with an optimal training to testing ratio of 4:1 (from previous section), using three different random seeds, and exclusively use patients of European descent for the entire training set. For comparison, we then include patients from all ethnic backgrounds in the training set. We use the same three seeds for the optimal train: test ratio of 4:1. To keep the training set size comparable, we reduce the number of white patients by 72% for each disease. In both scenarios, we limit the testing dataset to include only patients whose diagnostic diseases are already present in the training set. The results for the top-1 and top-10 accuracies are shown in **Fig. 2C**, labeled as "Top-1 (European)" and "Top-10 (European)," respectively.

Complete results, including Top-1, Top-10, Top-50, and Top-100 mean percentage accuracies for each individual ethnicity, along with their corresponding standard deviations, are documented in **Tables S4-S9**.

#### Training and testing of GestaltMML

We fine-tuned our training set on ViLT (see **Fig. 1**) and then tested on the various test sets. The testing data may include both data from GMDB and external validation data. The Vision-and-Language Transformer (ViLT)<sup>41</sup> utilized transformer encoder as modality interaction module and was pretrained on four datasets: Microsoft COCO (MSCOCO)<sup>76</sup>, Visual Genome (VG)<sup>77</sup>, SBU Captions (SBU)<sup>78</sup>, and Google Conceptual Captions (GCC)<sup>79</sup>. The statistics of these four datasets were reported in Table 1 of the original paper<sup>41</sup>.

### UMAP Clustering Analysis of GestaltMML.

In **Fig. 4**, showcasing the UMAP clustering outcomes, we employed the logit values derived from the GestaltMML's penultimate layer (immediately preceding the final softmax layer), resulting in a matrix of dimensions  $n \times 528$ , where "n" represents the total number of test samples. For the UMAP fitting on this  $n \times 528$  matrix, composed of stacked logit values, we configured the parameters as follows:  $n\_neighbors = 7$ ,  $min\_dist = 0.1$ , and  $n\_components = 2$ .

For UMAP clustering of external validation data (**Fig. 4A** and **Fig. 4C**), we used GestaltMML trained with an optimal train-to-test ratio of 9:1. For the comparison between NAA10 and NAA15 (**Fig. 4B**), to ensure a sufficient number of samples in the test set, we selected GestaltMML with an optimal train-to-test ratio of 4:1.

## **ACKNOWLEDGEMENTS**

We thank the GestaltMatcher Database (GMDB) which provides a collection of curated medical photography of genetic syndromes for training the multimodal model used in the current study. We thank patients and their families for contributing facial photos and phenotype descriptions to enable the establishment of computational models. We thank Mian Umair Ahsan and the IDDRC Biostatistics and Data Science core (HD105354) for consultation on machine-learning. This project is supported by NIH grant HG012655, HG013031, the CHOP Research Institute, the Lorenzo "Turtle" Sartini, Jr. Endowed Chair in Beckwith-Wiedemann Syndrome Research, and the Victoria Fertitta Fund through the Lorenzo "Turtle" Sartini Jr. Endowed Chair in Beckwith-Wiedemann Syndrome Research. Collection of photos for NAA10-related neurodevelopmental syndrome was supported by New York State Office for People with Developmental Disabilities (OPWDD) and NIH NIGMS R35GM133408. We would like to thank Steven Klein and Andrew George for helping to organize the photos and the phenotypic data for the patients with Beckwith-Wiedemann syndrome and Sotos syndrome.

## **DECLARATIONS**

### ***Availability of data and materials***

The GMDB (v1.0.9) database used in the current study can be obtained from <https://db.gestaltmatcher.org/>. All the software tools and computational workflow (as Jupyter Notebook) can be found at <https://github.com/WGLab/GestaltMML>. This study did not generate any new material.

### ***Code availability***

All the code can be accessed publicly on the following GitHub repository: <https://github.com/WGLab/GestaltMML>.

### ***Competing interests***

The authors declare no competing interests.



## REFERENCES

1. Ferreira CR. The burden of rare diseases. *American journal of medical genetics Part A*. 2019;179(6):885-892.
2. The Lancet Diabetes E. Rare diseases: individually rare, collectively common. *Lancet Diabetes Endocrinol*. 2023;11(3):139.
3. Genetic and Rare Disease Information Center. <https://rarediseases.info.nih.gov/about>. 2023.
4. Weinreich SS, Mangon R, Sikkens JJ, Teeuw ME, Cornel MC. Orphanet: a European database for rare diseases. *Ned Tijdschr Geneesk*. 2008;152(9):518-519.
5. McKusick VA. Mendelian Inheritance in Man and its online version, OMIM. *Am J Hum Genet*. 2007;80(4):588-604.
6. Bauskis A, Strange C, Molster C, Fisher C. The diagnostic odyssey: insights from parents of children living with an undiagnosed condition. *Orphanet J Rare Dis*. 2022;17(1):233.
7. Wu AC, McMahon P, Lu C. Ending the Diagnostic Odyssey-Is Whole-Genome Sequencing the Answer? *JAMA Pediatr*. 2020;174(9):821-822.
8. Sawyer SL, Hartley T, Dymant DA, et al. Utility of whole-exome sequencing for those near the end of the diagnostic odyssey: time to address gaps in care. *Clin Genet*. 2016;89(3):275-284.
9. Miller D. The diagnostic odyssey: our family's story. *Am J Hum Genet*. 2021;108(2):217-218.
10. Roosenboom J, Hens G, Mattern BC, Shriver MD, Claes P. Exploring the Underlying Genetics of Craniofacial Morphology through Various Sources of Knowledge. *Biomed Res Int*. 2016;2016:3054578.
11. Gurovich Y, Hanani Y, Bar O, et al. Identifying facial phenotypes of genetic disorders using deep learning. *Nature medicine*. 2019;25(1):60-64.
12. Hsieh T-C, Bar-Haim A, Moosa S, et al. GestaltMatcher facilitates rare disease matching using facial phenotype descriptors. *Nature genetics*. 2022;54(3):349-357.
13. Hustinx A, Hellmann F, Sümer Ö, et al. Improving Deep Facial Phenotyping for Ultra-rare Disorder Verification Using Model Ensembles. Paper presented at: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision2023.
14. Yi D, Lei Z, Liao S, Li SZ. Learning face representation from scratch. *arXiv preprint arXiv:14117923*. 2014.
15. Duta IC, Liu L, Zhu F, Shao L. Improved residual networks for image and video recognition. Paper presented at: 2020 25th International Conference on Pattern Recognition (ICPR)2021.
16. Deng J, Guo J, Xue N, Zafeiriou S. Arcface: Additive angular margin loss for deep face recognition. Paper presented at: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition2019.

17. Cao Q, Shen L, Xie W, Parkhi OM, Zisserman A. Vggface2: A dataset for recognising faces across pose and age. Paper presented at: 2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018)2018.
18. An X, Zhu X, Xiao Y, et al. Partialfc: Training 10 million identities on a single machine. arXivpreprint. *arXiv preprint arXiv:201005222*. 2020;3(4).
19. Şıklar Z, Berberoğlu M. Syndromic disorders with short stature. *Journal of clinical research in pediatric endocrinology*. 2014;6(1):1.
20. Wu L-T, Woody GE, Yang C, Pan J-J, Blazer DG. Racial/ethnic variations in substance-related disorders among adolescents in the United States. *Archives of General Psychiatry*. 2011;68(11):1176-1185.
21. Andresen EM, Brownson RC. Disability and health status: Ethnic differences among women in the United States. *Journal of Epidemiology & Community Health*. 2000;54(3):200-206.
22. Franks P, Gold MR, Fiscella K. Sociodemographics, self-rated health, and mortality in the US. *Social science & medicine*. 2003;56(12):2505-2514.
23. Fuller-Thomson E, Nuru-Jeter A, Minkler M, Guralnik JM. Black—White disparities in disability among older Americans: Further untangling the role of race and socioeconomic status. *Journal of aging and health*. 2009;21(5):677-698.
24. Chan ML, Eng CW, Gilsanz P, et al. Prevalence of Instrumental Activities of Daily Living Difficulties and Associated Cognitive Predictors Across Racial/Ethnic Groups: Findings From the KHANDLE Study. *The Journals of Gerontology: Series B*. 2022;77(5):885-894.
25. Nuru-Jeter AM, Thorpe Jr RJ, Fuller-Thomson E. Black-white differences in self-reported disability outcomes in the US: early childhood to older adulthood. *Public Health Reports*. 2011;126(6):834-843.
26. Walker JL, Harrison TC, Brown A, Thorpe Jr RJ, Szanton SL. Factors associated with disability among middle-aged and older African American women with osteoarthritis. *Disability and health journal*. 2016;9(3):510-517.
27. Krause JS, Broderick LE, Saladin L, Broyles J. Racial disparities in health outcomes after spinal cord injury: mediating effects of education and income. *The journal of spinal cord medicine*. 2006;29(1):17-25.
28. Green CA. Race, ethnicity, and Social Security retirement age in the US. *Feminist Economics*. 2005;11(2):117-143.
29. Hsieh T-C, Mensah MA, Pantel JT, et al. PEDIA: prioritization of exome data by image analysis. *Genetics in Medicine*. 2019;21(12):2807-2814.
30. Dingemans AJ, Hinne M, Truijien KM, et al. PhenoScore quantifies phenotypic variation for rare genetic diseases by combining facial analysis with other clinical features using a machine-learning framework. *Nature Genetics*. 2023;55(9):1598-1607.
31. Gilbert S, Harvey H, Melvin T, Vollebregt E, Wicks P. Large language model AI chatbots require approval as medical devices. *Nature Medicine*. 2023;29(10):2396-2398.
32. Vaswani A, Shazeer N, Parmar N, et al. Attention is all you need. *Advances in neural information processing systems*. 2017;30.

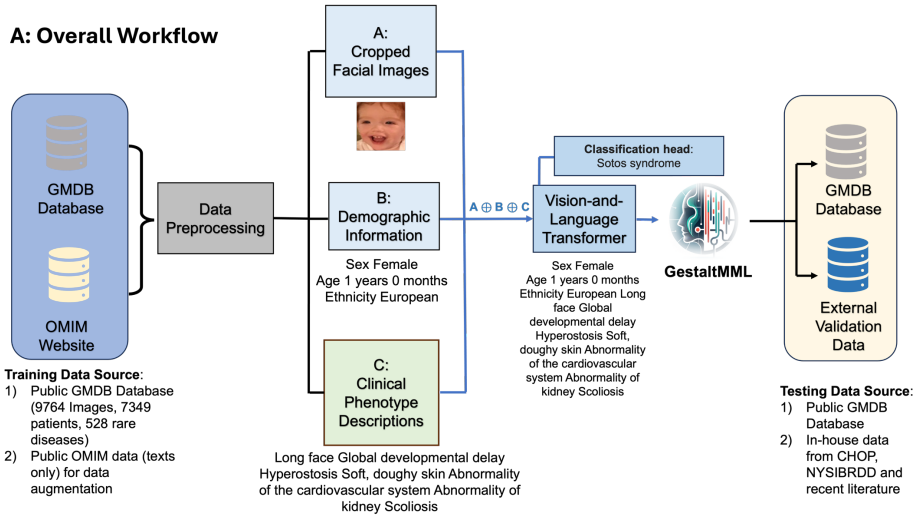
33. Jiang K, Lu X. Natural language processing and its applications in machine translation: a diachronic review. Paper presented at: 2020 IEEE 3rd International Conference of Safe Production and Informatization (IICSPI)2020.
34. Li J, Tang T, Zhao WX, Nie J-Y, Wen J-R. Pretrained language models for text generation: A survey. *arXiv preprint arXiv:220105273*. 2022.
35. Wankhade M, Rao ACS, Kulkarni C. A survey on sentiment analysis methods, applications, and challenges. *Artificial Intelligence Review*. 2022;55(7):5731-5780.
36. Yadav V, Bethard S. A survey on recent advances in named entity recognition from deep learning models. *arXiv preprint arXiv:191011470*. 2019.
37. Dosovitskiy A, Beyer L, Kolesnikov A, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:201011929*. 2020.
38. Carion N, Massa F, Synnaeve G, Usunier N, Kirillov A, Zagoruyko S. End-to-end object detection with transformers. Paper presented at: European conference on computer vision2020.
39. Strudel R, Garcia R, Laptev I, Schmid C. Segmenter: Transformer for semantic segmentation. Paper presented at: Proceedings of the IEEE/CVF international conference on computer vision2021.
40. Wang Y, Xu J, Sun Y. End-to-end transformer based model for image captioning. Paper presented at: Proceedings of the AAAI Conference on Artificial Intelligence2022.
41. Kim W, Son B, Kim I. Vilt: Vision-and-language transformer without convolution or region supervision. Paper presented at: International Conference on Machine Learning2021.
42. Radford A, Kim JW, Hallacy C, et al. Learning transferable visual models from natural language supervision. Paper presented at: International conference on machine learning2021.
43. Li LH, Yatskar M, Yin D, Hsieh C-J, Chang K-W. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:190803557*. 2019.
44. Li J, Selvaraju R, Gotmare A, Joty S, Xiong C, Hoi SCH. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in neural information processing systems*. 2021;34:9694-9705.
45. Gemini team. Gemini: A Family of Highly Capable Multimodal Models. *Google DeepMind* 2023; [https://storage.googleapis.com/deepmind-media/gemini/gemini\\_1\\_report.pdf](https://storage.googleapis.com/deepmind-media/gemini/gemini_1_report.pdf). Available at.
46. Lesmann H, Lyon GJ, Caro P, et al. GestaltMatcher Database-a FAIR database for medical imaging data of rare disorders. *MedRxiv*. 2023.
47. Edmondson AC, Kalish JM. Overgrowth syndromes. *Journal of pediatric genetics*. 2015:136-143.
48. Duffy KA, Cielo CM, Cohen JL, et al. Characterization of the Beckwith-Wiedemann spectrum: Diagnosis and management. Paper presented at: American Journal of Medical Genetics Part C: Seminars in Medical Genetics2019.

49. Brioude F, Kalish JM, Mussa A, et al. Clinical and molecular diagnosis, screening and management of Beckwith–Wiedemann syndrome: an international consensus statement. *Nature Reviews Endocrinology*. 2018;14(4):229-249.
50. Brioude F, Toutain A, Giabicani E, Cottureau E, Cormier-Daire V, Netchine I. Overgrowth syndromes—clinical and molecular aspects and tumour risk. *Nature Reviews Endocrinology*. 2019;15(5):299-311.
51. Duffy KA, Sajorda BJ, Yu AC, et al. Beckwith–Wiedemann syndrome in diverse populations. *American Journal of Medical Genetics Part A*. 2019;179(4):525-533.
52. Shuman C, Kalish JM, Weksberg R. Beckwith-wiedemann syndrome. *GeneReviews*<sup>®</sup>[Internet]. 2023.
53. Eggermann T, Algar E, Lapunzina P, et al. Clinical utility gene card for: Beckwith-Wiedemann Syndrome. *Eur J Hum Genet*. 2014;22(3).
54. Carli D, Bertola C, Cardaropoli S, et al. Prenatal features in Beckwith-Wiedemann syndrome and indications for prenatal testing. *Journal of Medical Genetics*. 2021;58(12):842-849.
55. Wang KH, Kupa J, Duffy KA, Kalish JM. Diagnosis and management of Beckwith-Wiedemann syndrome. *Frontiers in pediatrics*. 2020;7:562.
56. Tatton-Brown K, Rahman N. Sotos syndrome. *European Journal of Human Genetics*. 2007;15(3):264-271.
57. Baujat G, Cormier-Daire V. Sotos syndrome. *Orphanet Journal of Rare Diseases*. 2007;2(1):1-6.
58. Cole T, Hughes H. Sotos syndrome: a study of the diagnostic criteria and natural history. *Journal of medical genetics*. 1994;31(1):20-32.
59. Leventopoulos G, Kitsiou-Tzeli S, Kritikos K, et al. A clinical study of Sotos syndrome patients with review of the literature. *Pediatric neurology*. 2009;40(5):357-364.
60. Martinez F, Marín-Reina P, Sanchis-Calvo A, et al. Novel mutations of NFIX gene causing Marshall-Smith syndrome or Sotos-like syndrome: one gene, two phenotypes. *Pediatric Research*. 2015;78(5):533-539.
61. Wu Y, Lyon GJ. NAA10-related syndrome. *Exp Mol Med*. 2018;50(7):1-10.
62. Rope AF, Wang K, Evjenth R, et al. Using VAAST to identify an X-linked disorder resulting in lethality in male infants due to N-terminal acetyltransferase deficiency. *Am J Hum Genet*. 2011;89(1):28-43.
63. Lyon GJ, Vedaie M, Beisheim T, et al. Expanding the phenotypic spectrum of NAA10-related neurodevelopmental syndrome and NAA15-related neurodevelopmental syndrome. *Eur J Hum Genet*. 2023;31(7):824-833.
64. Deardorff MA, Noon SE, Krantz ID. Cornelia de Lange syndrome. 2020.
65. Berney TP, Ireland M, Burn J. Behavioural phenotype of Cornelia de Lange syndrome. *Archives of Disease in Childhood*. 1999;81(4):333-336.
66. Kline AD, Krantz ID, Sommer A, et al. Cornelia de Lange syndrome: clinical review, diagnostic and scoring systems, and anticipatory guidance. *American journal of medical genetics part A*. 2007;143(12):1287-1296.

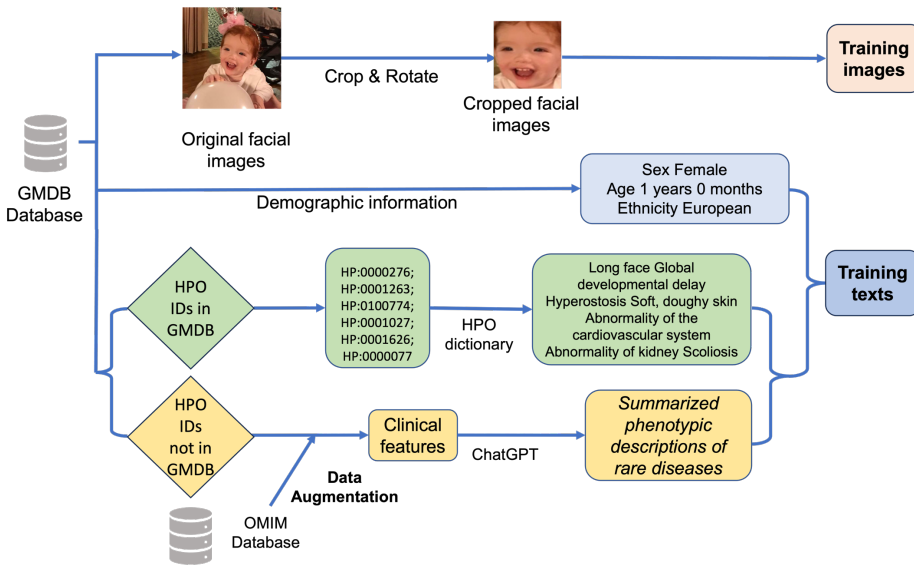
67. Morel Swols D, Foster J, Tekin M. KBG syndrome. *Orphanet journal of rare diseases*. 2017;12(1):1-7.
68. Low K, Ashraf T, Canham N, et al. Clinical and genetic aspects of KBG syndrome. *American journal of medical genetics Part A*. 2016;170(11):2835-2846.
69. Scarano E, Tassone M, Graziano C, et al. Novel mutations and unreported clinical features in KBG syndrome. *Molecular Syndromology*. 2019;10(3):130-138.
70. Gnazzo M, Lepri FR, Dentici ML, et al. KBG syndrome: common and uncommon clinical features based on 31 new patients. *American Journal of Medical Genetics Part A*. 2020;182(5):1073-1083.
71. Brancati F, D'Avanzo MG, Digilio MC, et al. KBG syndrome in a cohort of Italian patients. *American Journal of Medical Genetics Part A*. 2004;131(2):144-149.
72. Zeeberg BR, Qin H, Narasimhan S, et al. High-Throughput GoMiner, an 'industrial-strength' integrative gene ontology tool for interpretation of multiple-microarray experiments, with application to studies of Common Variable Immune Deficiency (CVID). *BMC Bioinformatics*. 2005;6:168.
73. Yang J, Liu C, Deng W, et al. Enhancing phenotype recognition in clinical notes using large language models: PhenoBCBERT and PhenoGPT. *Patterns*. 2024;5(1).
74. Zhong S, Storch KF, Lipan O, Kao MC, Weitz CJ, Wong WH. GoSurfer: a graphical interactive tool for comparative analysis of large gene sets in Gene Ontology space. *Appl Bioinformatics*. 2004;3(4):261-264.
75. OpenAI. ChatGPT. <https://openai.com/chatgpt>, 2023.
76. Lin T-Y, Maire M, Belongie S, et al. Microsoft coco: Common objects in context. Paper presented at: Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 132014.
77. Krishna R, Zhu Y, Groth O, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*. 2017;123:32-73.
78. Ordonez V, Kulkarni G, Berg T. Im2text: Describing images using 1 million captioned photographs. *Advances in neural information processing systems*. 2011;24.
79. Sharma P, Ding N, Goodman S, Soricut R. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. Paper presented at: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)2018.

# FIGURES

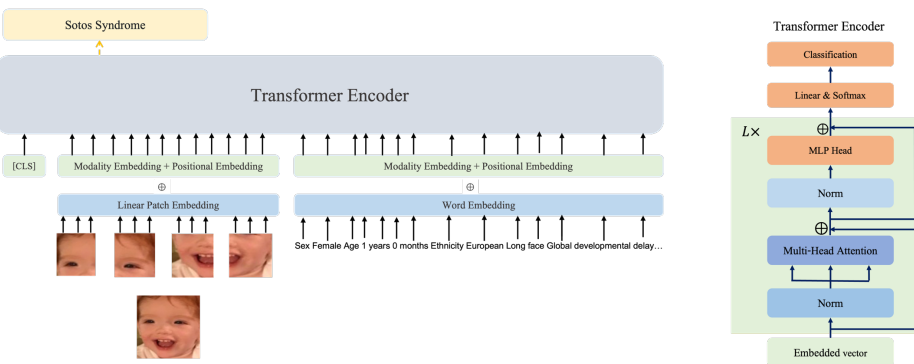
## A: Overall Workflow



## B: Data Preprocessing Pipeline

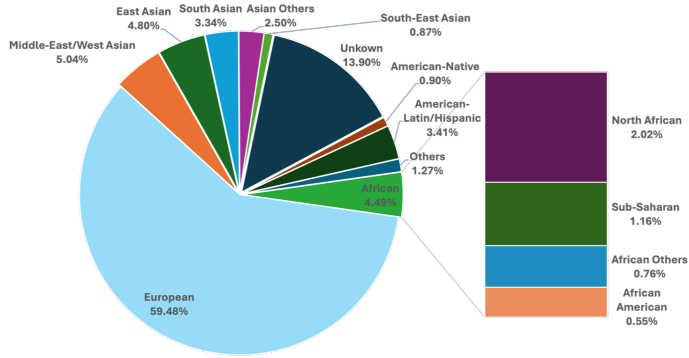


## C: Architectural Framework

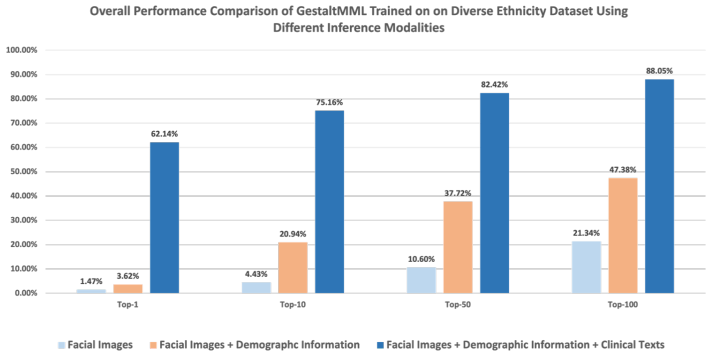
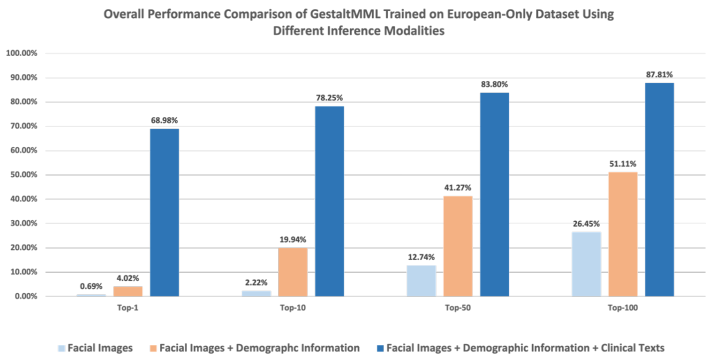


**Fig. 1. Overview of the GestaltMML. A: Illustration of the overall workflow of the project.** GestaltMML uses information from facial images after appropriate pre-processing, demographic information as well as description of clinical phenotypes on each disease from both GMDB (if available) and the OMIM database. GMDB: GestaltMatcher Database. OMIM: Online Mendelian Inheritance in Man. CHOP: Children’s Hospital of Philadelphia. NYSIBRDD: New York State Institute for Basic Research in Developmental Disabilities. **B: Data Preprocessing Pipeline of GestaltMML, using Sotos syndrome as an example.** The facial images in GMDB were cropped by “FaceCropper” to crop and rotate the size of 112 \* 112. The training texts can be divided into two categories: (1) Demographic information + HPO textual data, and (2) Demographic information + clinical features from OMIM database summarized by ChatGPT. **C: Architectural Framework of GestaltMML:** Based on the foundation of ViLT the structure of GestaltMML employs the Transformer encoder, capable of processing both textual and image inputs. Notice that this architecture closely resembles ViT, with the distinction that it solely accepts images as its input.

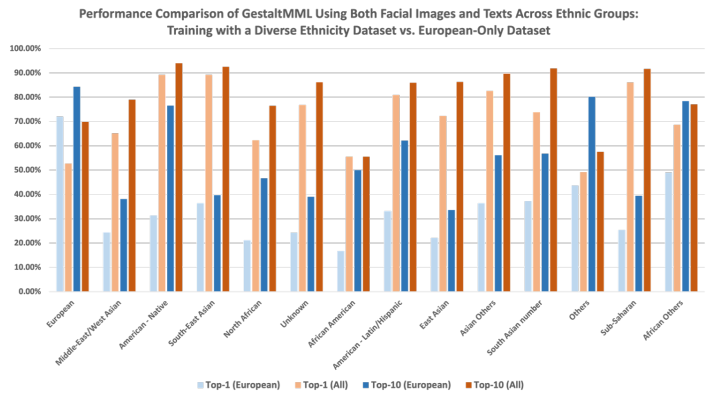
**A: Ethnicity distribution within GMDB (v1.0.9)**



**B: Performance comparison across different modalities**



**C: Performance comparison across different training sets**





**Fig. 2. A:** Ethnicity distribution in GMDB (v1.0.9). **B:** Overall accuracy by GestaltMML when using different modalities for inference. **C:** Comparative analysis of GestaltMML's effectiveness trained on patients of European descents only vs. patients of all ethnic backgrounds. Both set of experiments use the same set of training and testing sets. The specifics regarding the sizes of the training and testing datasets are documented in **Table S3**.

### A: Beckwith-Wiedemann Syndrome

Most-likely diseases among 528 total rare genetic diseases



Only use Cropped Facial Image



GestaltMML

.....  
**Rank 6: Beckwith-Wiedemann Syndrome** ✓  
 .....

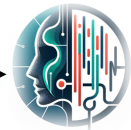


Cropped Facial Image

Sex male Age years months Ethnicity African  
 Macroglossia Hemimacroglossia Omphalocele  
 Hemihypertrophy Hemihypertrophy of lower limb  
 Hemihypertrophy of upper limb Large for gestational  
 age Linear earlobe crease Diagonal earlobe crease  
 Anterior creases of earlobe Postauricular pit  
 Preauricular pit Posterior helix pit Supraauricular pit  
 Hypoglycemia Neonatal hypoglycemia

Demographic Information

Clinical Phenotype Descriptions



GestaltMML

.....  
**Rank 1: Beckwith-Wiedemann Syndrome** ✓  
 .....

### B: Sotos Syndrome

Most-likely diseases among 528 total rare genetic diseases

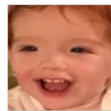


Only use Cropped Facial Image



GestaltMML

.....  
**Rank 288: Sotos Syndrome** ✓  
 .....

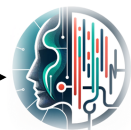


Cropped Facial Image

Sex Female Age 1 years 0 months Ethnicity European  
 Long face Global developmental delay Hyperostosis Soft,  
 doughy skin Abnormality of the cardiovascular system  
 Abnormality of the kidney Scoliosis

Demographic Information

Clinical Phenotype Descriptions

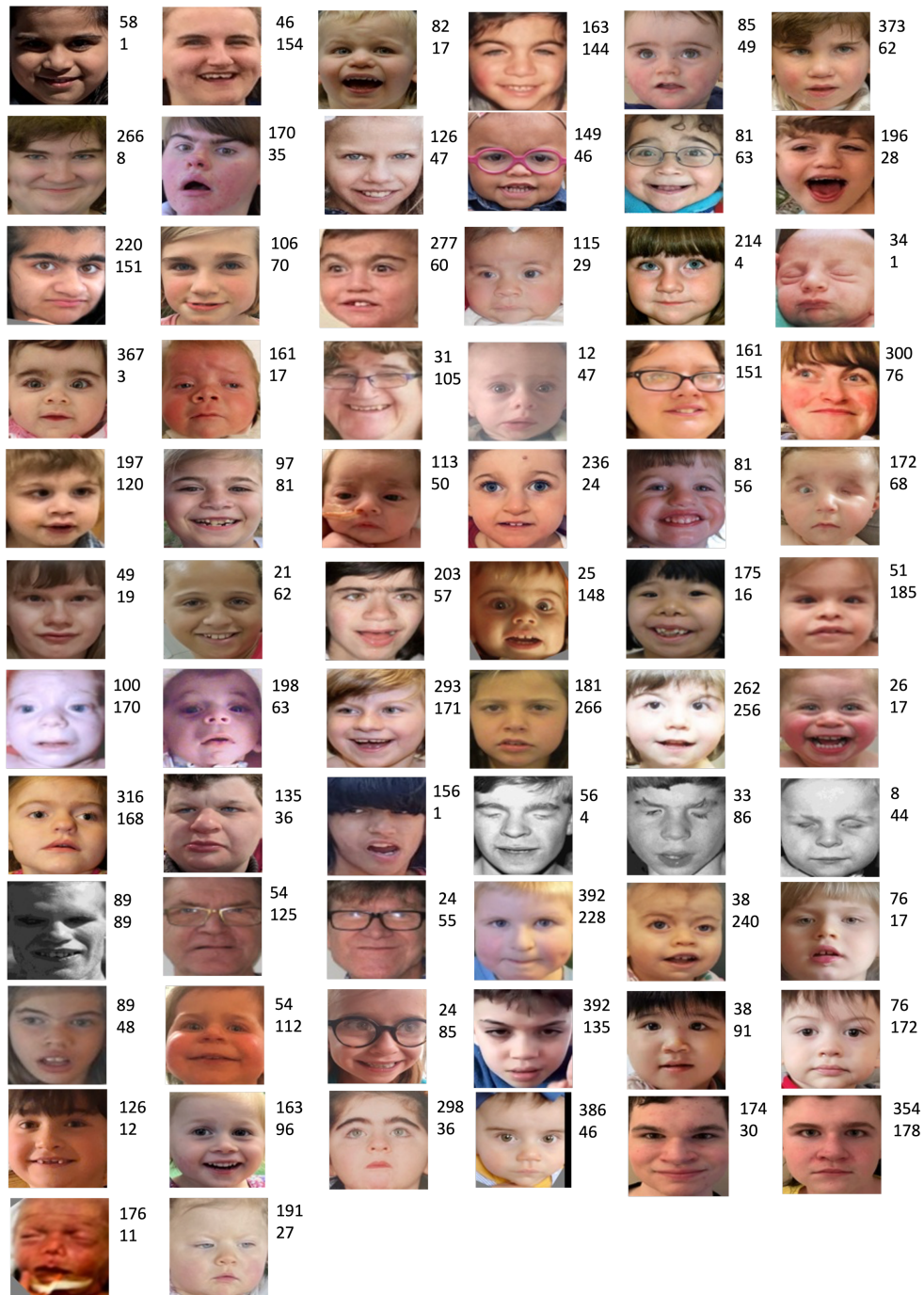


GestaltMML

.....  
**Rank 1: Sotos Syndrome** ✓  
 .....

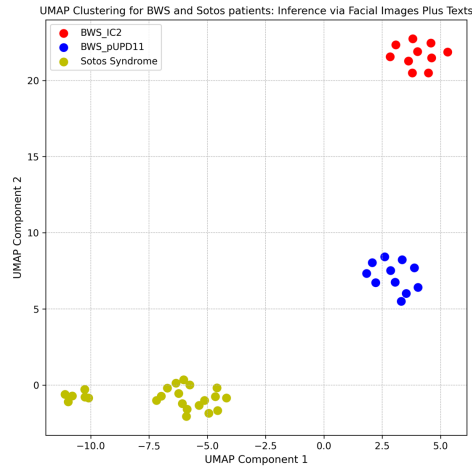
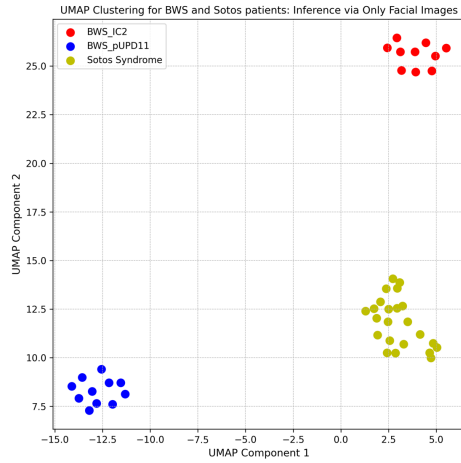
**Fig. 3. Illustration demonstrating that multimodal inference, which combines facial images, demographic details, and clinical phenotype descriptions, is more effective than using facial images alone in diagnosing patients with BWS (A) and Sotos Syndrome (B). The patient data for this study were sourced from CHOP with the appropriate consent.**

**Upper number:** rank of true label predicted by GestaltMML (v1.0.3) within 449 rare diseases using **facial images only**.  
**Lower number:** rank of true label predicted by GestaltMML (v1.0.3) within 449 rare diseases using **facial images and texts**.

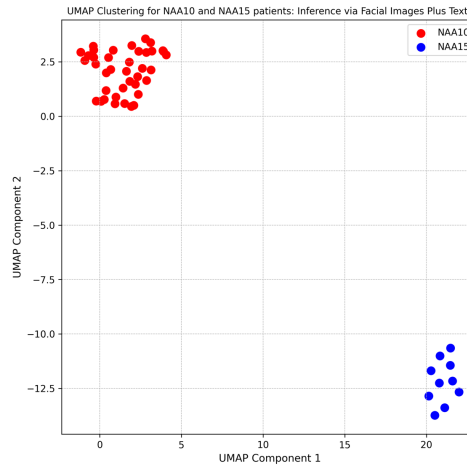
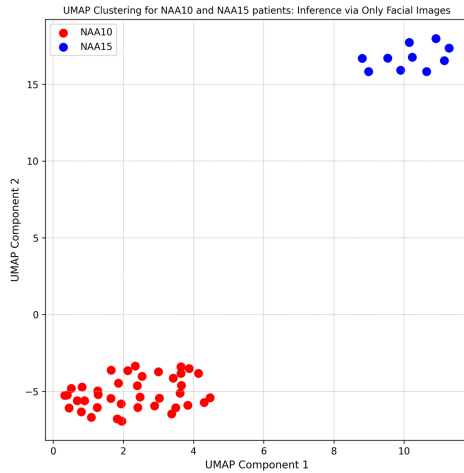


**Figure 4. Rank of true label among 449 total disease labels for 68 NAA10 patients (with proper consent) from NYSIBRDD and recently published literature, as predicted by GestaltMML (trained on GMDB v1.0.3. with optimal train: test ratio of 4:1). The number on the upper and lower level indicate the ranking achieved using only facial images and using combined information (facial images, demographic data and clinical phenotype descriptions), respectively.**

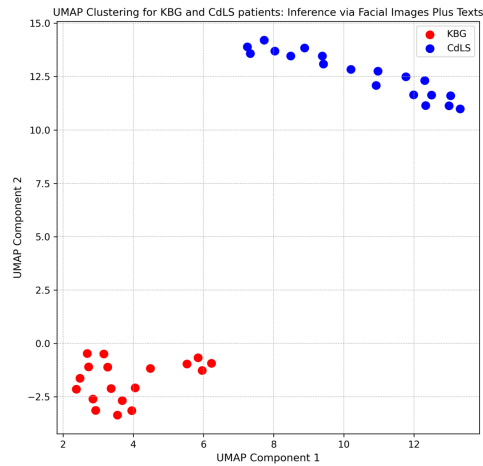
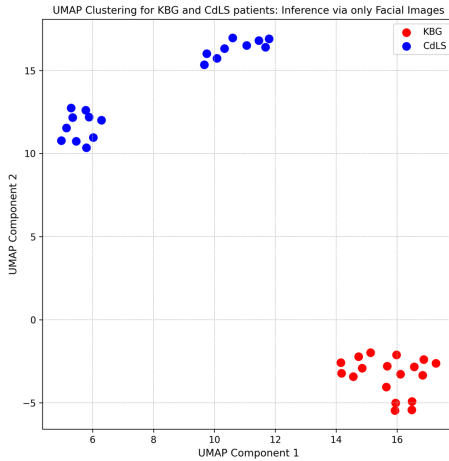
### A: BWS vs. Sotos



### B: NAA10 vs. NAA15



### C: KBG vs. CdLS



**Fig. 5. UMAP Clustering Analysis Across Three Comparative Sets. A: BWS Patient Cohorts (IC2 and pUPD11) Alongside Sotos Syndrome Patients B: NAA10 and NAA15 patients within GMDB (v1.0.9) C: KBG Syndrome and CdLS patients. For each comparison: On the left, inference is performed using only facial images; on the right, a multimodal approach combines facial images with textual data.**

## TABLES

**Table 1.** Overview of the GMDB (v1.0.9) dataset, including the GMDB-frequent subset and the GMDB-rare subset. The GMDB-frequent contains disorders with more than ( $>$ ) 6 patients, while the GMDB-rare contains disorders with less or equal to ( $\leq$ ) 6 patients.

Dataset	# of images	# of patients	# of disorders	# of images with clinical HPO texts
<b>GMDB-frequent</b>	8547	6376	244	3962
<b>GMDB-rare</b>	1217	973	284	470
<b>Total</b>	9764	7349	528	4432

**Table 2. Performance of GestaltMML, GestaltViT and GestaltLT and their comparisons with that of baseline Ensembled image model** in <sup>13</sup>: All the models have undergone fine-tuning on the GMDB (v1.0.9). The models with underline indicates that we *exclusively* tested them on a subset of the total test set containing the available test modalities.

Model	GMDB-frequent		GMDB-rare	
	Top-1	Top-10	Top-1	Top-10
<b><u>GestaltMML</u></b>	50.14%	75.50%	24.14%	41.38%
<b><u>GestaltLT</u></b>	46.40%	74.93%	24.16%	41.38%
<b>GestaltViT</b>	18.16%	44.67%	6.97%	18.12%
<b>Ensembled image model</b>	43.02%	72.44%	19.77%	38.57%

**Table 3.** Results of GestaltMML on patients diagnosed with Beckwith-Wiedemann Syndrome (BWS), Sotos Syndrome and Cornelia de Lange Syndrome (CdLS) at the Children’s Hospital of Philadelphia, and KBG syndrome at the New York State Institute for Basic Research in Developmental Disabilities and recently published literature.

<b>Beckwith-Wiedemann Syndrome (BWS)</b>								
Model	Cohort (subset)	Testing Modalities	Percentage Accuracy				Sample Size (outside validation)	# of images in GMDB (v1.0.9)
			Top-1	Top-10	Top-50	Top-100		
GestaltMML	IC2	Images + Texts	100%	100%	100%	100%	10	26
GestaltMML	IC2	Texts only	90%	100%	100%	100%	10	26
GestaltMML	IC2	Images only	10%	40%	70%	90%	10	26
Ensembled image model	IC2	Images only	20%	60%	90%	100%	10	26
GestaltMML	pUPD 11	Images + Texts	100%	100%	100%	100%	11	26
GestaltMML	pUPD 11	Texts only	100%	100%	100%	100%	11	26
GestaltMML	pUPD 11	Images only	0%	27.27%	45.45%	54.54%	11	26
Ensembled image model	pUPD 11	Images only	8.33%	66.67%	91.67%	100%	11	26
<b>Sotos Syndrome</b>								
Model	Testing Modalities	Percentage Accuracy				Sample Size (outside validation)	# of images in GMDB (v1.0.9)	
		Top-1	Top-10	Top-50	Top-100			
GestaltMML	Images + Texts	73.91%	86.96%	95.65%	95.65%	23	126	
GestaltMML	Texts only	69.57%	82.61%	95.65%	95.65%	23	126	
GestaltMML	Images only	0%	4.34%	17.39%	30.43%	23	126	

Ensembled image model	Images only	41.93%	61.29%	93.55%	95.16%	23	126
<b>NAA10-related Neurodevelopmental Syndrome</b>							
Model	Testing Modalities	Percentage Accuracy				Sample Size (outside validation)	# of images in GMDB (v1.0.3)
		Top-1	Top-10	Top-50	Top-100		
GestaltMML	Images + Texts	4.41%	32.35%	66.18 %	82.35%	68	15
GestaltMML	Texts only	5.88%	38.23%	72.06%	82.35%	68	15
GestaltMML	Images only	0.00%	1.47%	23.53%	41.18%	68	15
Ensembled image model	Images only	7.35%	16.18%	44.12%	66.18%	68	15

<b>Cornelia de Lange Syndrome (CdLS)</b>							
Model	Testing Modalities	Percentage Accuracy				Sample Size (outside validation)	# of images in GMDB (v1.0.9)
		Top-1	Top-10	Top-50	Top-100		
GestaltMML	Images + Texts	21.05%	73.68%	84.21%	94.74%	19	382
GestaltMML	Texts only	0.00%	47.36%	84.21%	89.47%	19	382
GestaltMML	Images only	52.63%	68.42%	84.21%	94.74%	19	382
Ensembled image model	Images only	76.67%	86.67%	100%	100%	19	382

KBG Syndrome							
Model	Testing Modalities	Percentage Accuracy				Sample Size (outside validation)	# of images in GMDB (v1.0.9)
		Top-1	Top-10	Top-50	Top-100		
GestaltMML	Images + Texts	44.44%	83.33%	88.89%	88.89%	18	167
GestaltMML	Texts only	38.89%	72.22%	88.89%	88.89%	18	167
GestaltMML	Images only	55.56%	66.67%	83.33%	88.89%	18	167
Ensembled image model	Images only	94.44%	94.44%	100%	100%	18	167