

Supplementary Information for “Phenotype Driven Molecular Genetic Test Recommendation for Diagnosing Pediatric Rare Disorders”

Fangyi Chen¹, Priyanka Ahimaz^{2,3}, Quan M. Nguyen^{4,5}, Rachel Lewis², Wendy K. Chung⁶, Casey Ta¹, Katherine M. Szigety⁷, Sarah E. Sheppard⁷, Ian M. Campbell⁷, Kai Wang⁴, Chunhua Weng^{1,*}, Cong Liu^{1,*}

¹Department of Biomedical Informatics, Columbia University, New York, NY, USA;

²Department of Pediatrics, Columbia University, New York, NY, USA;

³Institute of Genomic Medicine, Columbia University, New York, NY, USA;

⁴Raymond G. Perelman Center for Cellular and Molecular Therapeutics, Children’s Hospital of Philadelphia, Philadelphia, PA, USA;

⁵Department of Bioengineering, University of Pennsylvania, Philadelphia, PA, USA;

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

⁷Division of Human Genetics, Department of Pediatrics, Children’s Hospital of Philadelphia, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA;

(*: equal-contribution senior corresponding authors; emails: Dr. Cong Liu -

c13720@cumc.columbia.edu; Dr. Chunhua Weng - cw2384@cumc.columbia.edu)

Table of Contents

<i>Supplementary Table 1. Example of patient clinical summaries</i>	3
<i>Supplementary Table 2. Clinical summary of patients with test order adjusted to ES/GS.</i>	4
<i>Supplementary Table 4. Lists of keywords used to filter genetic-related notes, and search for molecular tests ordered by clinicians, as well as exceptions words for panel group were being noted.</i>	6
<i>Supplementary Data</i>	7
Supplementary Data 1. Categories of gene panels in Columbia University Irving Medical Center (CUIMC) initial genetic cohort.	7
Supplementary Data 2. Experimental results of various combinations of models, sorted by average areas under precision-recall curve (AUPRC) in descending order. Performance metrics were averaged and displayed below.....	7
Supplementary Data 3. Feature Importance of HPO phenotypic abnormalities calculated by Gini impurity, along with OLS estimated coefficients.....	7
Supplementary Data 4. Feature importance pf phecodes calculated by Gini impurity, along with phecode sets and total sum of feature importances by corresponding phecodes.....	7
Supplementary Data 5. Model Performance within subgroups characterized by demographic characteristics.....	7
Supplementary Data 6. Lists of OMOP concept IDs used to identify the "genetic" cohort from the structured OMOP database.	7

Supplementary Table 1. Example of patient clinical summaries. The clinical summaries entail key phenotype indicators that genetic experts believe are relevant to the genetic disorders. It is important to note that the phenotype summary was manually recorded in the research database, separate from the original EHR data (used for model training), and was not directly utilized as features for model training.

MRN	Primary Indication	Test	Test Date
xxxxxxx	ACL tear	Gene panel (Ehlers Danlos Syndrome panel)	xx/xx/xxxx
xxxxxxx	Marinesco-Sjögren syndrome and a history of catatonia	Whole Exome	xx/xx/xxxx
xxxxxxx	pulmonary hypertension, severe hydronephrosis, and posterior urethral valves, status post repair	Whole Genome	xx/xx/xxxx

Supplementary Table 2. Clinical summary of patients with test order adjusted to ES/GS.

Primary Indication	Counts
Seizures	59
Autism spectrum disorder	52
Congenital heart defect	13
Developmental delay	10
Multiple birth defects	5

Supplementary Table 3. Features used to train the classification models.

Categories	Feature Names & Description	Input Feature Dimension
Clinical Features (Structured Data)	- <i>Freq_phecodes</i> : frequency of each phecodes/phenotype	1,225
	- <i>Sum_phecodes</i> : total number of unique phecodes/phenotypes	1
	- <i>Freq_HPO</i> : frequency of each HPO-based organ systems of phenotypic abnormality	23
Clinical Features (Unstructured Data)	- <i>Freq_phecodes_notes</i> : frequency of each phecodes/phenotypes derived from clinical narratives	418
	- <i>Sum_phecodes_notes</i> : total number of unique phecodes/phenotypes derived from clinical narratives	1
	- <i>Num_notes</i> : cumulative sum of notes	1
Demographics Characteristics	- <i>Age</i>	1
	- Sex assigned at birth time	1
	- <i>Race self-reported by patients</i>	1

Supplementary Table 4. Lists of keywords used to filter genetic-related notes, and search for molecular tests ordered by clinicians, as well as exceptions words for panel group were being noted.

Keywords in Note Titles	Keywords in label determination	Exceptions in Panel Group*
"genetic", "letter", "progress", "visit", "progress"	"WES", "WGS", "exome", "genomic", "panel"	blood', 'screen', 'screening', 'viral', 'virus', 'pcr', 'metabolic', 'hepatic', 'lipid', 'tcell', 't cell', 't-cell', 'iron', 'respiratory', "pathogen", "feeding", "liver", "thyroid", "immunoglobulin", "allergy", "allergen", "celiac", 'antigen', "hepatitis", 'vitamin', "chemistry"

Supplementary Data

Please refer Supplementary Data to excel sheet.

Supplementary Data 1. Categories of gene panels in Columbia University Irving Medical Center (CUIMC) initial genetic cohort.

Supplementary Data 2. Experimental results of various combinations of models, sorted by average areas under precision-recall curve (AUPRC) in descending order. Performance metrics were averaged and displayed below.

Supplementary Data 3. Feature Importance of HPO phenotypic abnormalities calculated by Gini impurity, along with OLS estimated coefficients.

Supplementary Data 4. Feature importance of phecodes calculated by Gini impurity, along with phecode sets and total sum of feature importances by corresponding phecodes.

Supplementary Data 5. Model Performance within subgroups characterized by demographic characteristics.

Supplementary Data 6. Lists of OMOP concept IDs used to identify the "genetic" cohort from the structured OMOP database.