

## RESEARCH ARTICLE

## HUMAN GENOMICS

# Phenotype risk scores identify patients with unrecognized Mendelian disease patterns

Lisa Bastarache,<sup>1</sup> Jacob J. Hughey,<sup>1</sup> Scott Hebbbring,<sup>2</sup> Joy Marlo,<sup>1</sup> Wanke Zhao,<sup>3</sup> Wanting T. Ho,<sup>3</sup> Sara L. Van Driest,<sup>4,5</sup> Tracy L. McGregor,<sup>5</sup> Jonathan D. Mosley,<sup>4</sup> Quinn S. Wells,<sup>4,6</sup> Michael Temple,<sup>1</sup> Andrea H. Ramirez,<sup>4</sup> Robert Carroll,<sup>1</sup> Travis Osterman,<sup>1,4</sup> Todd Edwards,<sup>4</sup> Douglas Ruderfer,<sup>4</sup> Digna R. Velez Edwards,<sup>7</sup> Rizwan Hamid,<sup>5</sup> Joy Cogan,<sup>5</sup> Andrew Glazer,<sup>4</sup> Wei-Qi Wei,<sup>1</sup> QiPing Feng,<sup>6</sup> Murray Brilliant,<sup>2</sup> Zhizhuang J. Zhao,<sup>3</sup> Nancy J. Cox,<sup>4</sup> Dan M. Roden,<sup>1,4,6</sup> Joshua C. Denny<sup>1,4,\*</sup>

Genetic association studies often examine features independently, potentially missing subpopulations with multiple phenotypes that share a single cause. We describe an approach that aggregates phenotypes on the basis of patterns described by Mendelian diseases. We mapped the clinical features of 1204 Mendelian diseases into phenotypes captured from the electronic health record (EHR) and summarized this evidence as phenotype risk scores (PheRSs). In an initial validation, PheRS distinguished cases and controls of five Mendelian diseases. Applying PheRS to 21,701 genotyped individuals uncovered 18 associations between rare variants and phenotypes consistent with Mendelian diseases. In 16 patients, the rare genetic variants were associated with severe outcomes such as organ transplants. PheRS can augment rare-variant interpretation and may identify subsets of patients with distinct genetic causes for common diseases.

Classically, Mendelian diseases are thought to be rare, caused by variants with large effect sizes, and associated with considerable morbidity and mortality. Many are characterized by a range of clinical phenotypes, often affecting multiple organ systems. Several lines of evidence suggest that genes known to cause Mendelian disease also harbor variants that contribute to complex disease (1). Studies have identified clinical overlap in patients codiagnosed with Mendelian and complex diseases (2), and single-nucleotide polymorphisms found in genome-wide association studies (GWASs) are enriched for Mendelian loci (3). A review of evidence from GWASs and whole-exome sequencing studies found an overlap between primary immunodeficiency genes and complex inflammatory diseases (4). Collectively, this evidence suggests that variants in Mendelian disease-causing genes may be an underrecognized contributor to complex disease.

Until very recently, the phenotypic effects of rare genetic variants were ascertained primarily through family-based studies of patients with distinctive and often severe phenotypes. Population-level techniques, such as GWASs and phenome-wide association studies (PheWASs) (5, 6), are not easily applied to rare variation because most studies are underpowered. Cohorts large enough to support GWASs of rare variants have only recently been assembled and have demonstrated the potential impact of rare variants on complex traits such as height, finding rare variants with effect sizes much greater than those of common variants (7).

Estimating the pathogenicity of rare variants remains a challenge and a barrier to use in the clinical setting (8). Many algorithms have been developed to predict variant pathogenicity (9–11), and consortia such as ClinGen (12) are aggregating knowledge to enable expert determinations. Resources such as ExAC (13) have helped refine variant interpretation. Some variants previously interpreted as pathogenic are too common in some populations to cause rare, life-threatening disorders (14), whereas others thought to be completely penetrant do not always cause disease (15). Initial studies suggest that electronic health records (EHRs) linked to genetic data may help drive genomic discovery and define clinical phenotypes associated with rare variants (16–18).

We have developed an approach that increases the power to detect rare-variant associations

by leveraging the phenotypic patterns of Mendelian diseases. By mapping the clinical manifestations of a Mendelian disease to phenotypes extracted from the EHR, we can compute a “phenotype risk score” (PheRS) that expresses the degree to which an individual’s symptoms overlap with a Mendelian disease. We defined a PheRS as a weighted aggregation of genetically related phenotypes, analogous to the genetic risk score approach for analyzing multiple variants against a single phenotype. PheRS was validated against clinically diagnosed cases and controls, and a genetic association study of PheRS profiles for 1204 Mendelian conditions identified both known and previously unknown associations with variants in target genes. The approach presents a method for measuring the phenotypic impact of rare variants and for identifying the heretofore underrecognized contribution of Mendelian disease genes to common medical conditions.

## Constructing a phenotype risk score

The Online Mendelian Inheritance in Man (OMIM) database provides clinical synopses for thousands of monogenic diseases (19) that have been annotated using the Human Phenotype Ontology (HPO) (20). We created a map from HPO terms to consolidated billing codes from the EHR called phecodes. Phecodes enable high-throughput ascertainment of EHR phenotypes and have been widely used to replicate known genetic associations and discover new ones (21–23). By mapping HPO terms to phecodes, we can express “phenotype syndromes” patterned after Mendelian diseases in OMIM in terms of clinical phenotypes that can be rapidly derived from the EHR. The PheRS for a given Mendelian disease is defined as the sum of clinical features observed in a given subject weighted by the log inverse prevalence of the feature.

## Validating PheRS

We compared the PheRSs of clinically diagnosed cases with matched controls for six Mendelian diseases. PheRS was a very strong predictor of case status for five of the diseases (Wilcoxon rank sum test;  $P = 8 \times 10^{-42}$  to  $5 \times 10^{-320}$ ) (Fig. 1, A and B). The exception was phenylketonuria ( $P = 0.28$ ), which effectively served as a negative control because newborn screening and dietary avoidance of phenylalanine essentially eliminates disease manifestations in affected individuals (24). The PheRS for each Mendelian disease demonstrated specificity for the target disease: The cases for different Mendelian diseases had similar PheRS distributions to those for controls (Fig. 1C). The lone exception was that the PheRS for hereditary hemochromatosis (HH) was significantly elevated for cystic fibrosis (CF) cases versus controls. However, even in this instance, CF cases had PheRSs that were three times as high for CF compared with HH. A review of controls with a PheRS greater than the 75th percentile identified one individual (PheRS > 99th percentile) who was diagnosed with HH in the 6 months after the case/control ascertainment. Thus, for this individual,

<sup>1</sup>Department of Biomedical Informatics, Vanderbilt University Medical Center, Nashville, TN, USA. <sup>2</sup>Center for Human Genetics, Marshfield Clinic Research Institute, Marshfield, WI, USA. <sup>3</sup>Department of Pathology, University of Oklahoma Health Sciences Center, Oklahoma City, OK, USA. <sup>4</sup>Department of Medicine, Vanderbilt University Medical Center, Nashville, TN, USA. <sup>5</sup>Department of Pediatrics, Vanderbilt University Medical Center, Nashville, TN, USA. <sup>6</sup>Department of Pharmacology, Vanderbilt University Medical Center, Nashville, TN, USA. <sup>7</sup>Department of Obstetrics and Gynecology, Vanderbilt University Medical Center, Nashville, TN, USA.

\*Corresponding author. Email: josh.denny@vanderbilt.edu

the PheRS suggested the diagnosis before it was made by providers.

PheRS identifies potentially pathogenic variants in Mendelian disease genes

We conducted an association analysis based on a cohort of 21,701 adults of European ancestry genotyped on the Exome BeadChip (table S1). In this cohort, we computed PheRSs for 1204 Mendelian diseases (1096 causative genes) for which we had sufficient genotype data. We tested for association between PheRSs and 6188 rare variants (minor allele frequency < 1%) using linear regression, assuming a dominant genetic model. We only tested the PheRS for a particular Mendelian disease against variants in the gene or genes known to cause that disease. We found 18 significant associations between rare variants and PheRSs (false discovery rate  $q < 0.1$ ; Table 1). All significant results had a positive beta coefficient, indicating that the variants were associated with an excess of Mendelian disease phenotypes. Four of the genes had an established dominant mode of inheritance, whereas the remaining 13 genes were known as exclusively or primarily recessive. Four were annotated in ClinVar as “pathogenic” or “likely pathogenic,” and the Human Gene Mutation Database (HGMD) provided evidence of pathogenicity for an additional three variants (25). The phenotypic impact of the remaining nine variants have not, to our knowledge, been previously described.

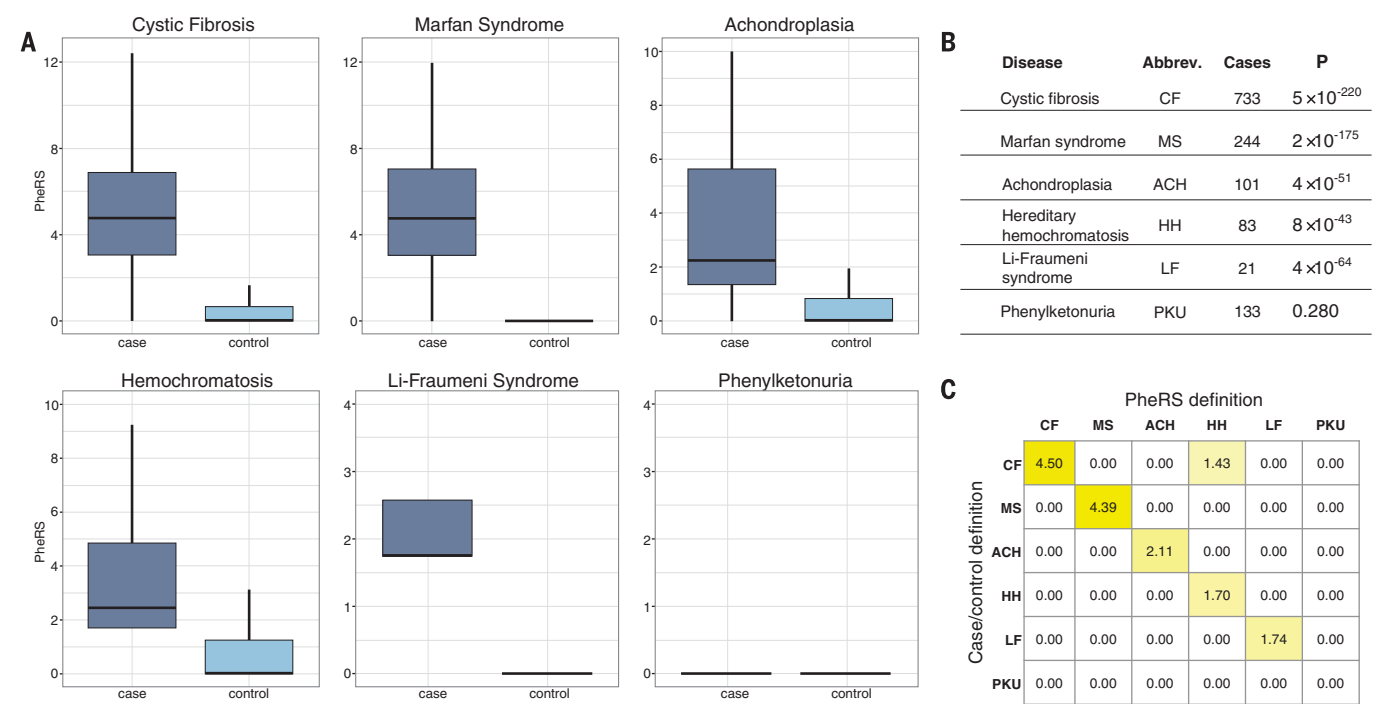
Clinical chart review revealed that eight of the 807 individuals with statistically significant var-

iants were diagnosed with the target Mendelian disease. Seven individuals with one of the two *CFTR* variants (p.G542\* and p.R553\*) were diagnosed with CF. Clinical genetic testing confirmed the variants called on the Exome BeadChip, including one homozygote for p.G542\*, and established compound heterozygosity with  $\Delta F508$  for five others (Fig. 2A and table S2). All individuals diagnosed with CF had a PheRS greater than four standard deviations from expected values. Additionally, the highest-scoring heterozygote for p.E168Q in *HFE* was diagnosed with HH on the basis of clinical findings. The diagnosis of HH was considered but never confirmed for another p.E168Q heterozygote who died of end-stage liver disease.

Although the majority of patients with significant variants were undiagnosed, these individuals had a high burden of severe end points related to the Mendelian diseases. Of the 40 heterozygotes for *HFE* p.E168Q, four had liver transplants (10 versus 1.2% in background;  $P = 2.1 \times 10^{-3}$ ; Fisher’s exact test). Individuals with variants in two genes associated with renal failure had elevated rates of kidney transplant: Five of 36 patients (14%) with *AGXT* p.A295T received transplants (another is awaiting transplant), as well as two of 15 patients (13%) with *DGKE* p.W322\* (versus 3% in background;  $P = 6.9 \times 10^{-3}$  and  $P = 0.088$ , respectively; Fisher’s exact test). Four of 69 *TG* p.G77S heterozygotes underwent thyroidectomies (6 versus 2% of noncarriers;  $P = 0.039$ ; Fisher’s exact test). These are end-stage phenotypes, potentially resulting from the effects of these variants, and

did not have an increased prevalence in other significant variant carriers (table S3). Additionally, we found that the population attributable fraction for the constituent PheRS phenotypes averaged 0.5%, with a maximum of 4.5%, suggesting that common diseases in adult populations may, in some cases, be attributed to variants in Mendelian genes (fig. S1).

An examination of PheWAS results using Fisher’s exact test for variants identified in the discovery analysis revealed that constituent phenotypes were often marginally significant ( $P < 0.05$ ), while not crossing the Bonferroni correction level for a single PheWAS. For *CFTR* p.G542\*, three features used in the PheRS for CF achieved marginal significance (bronchiectasis, disease of the pancreas, and chronic airway obstruction) (Fig. 2B). However, the constituent phenotypes for CF were only statistically significant when they were analyzed collectively as a PheRS. The association of p.G542\* with the PheRS for CF was similar to the association with the phenotype of CF itself (PheRS by linear regression,  $P = 3 \times 10^{-8}$ , versus CF diagnosis by Fisher’s exact test,  $P = 8 \times 10^{-7}$ ). Similarly, although individuals with variant p.W322\* in *DGKE* had an excess of nephrotic syndrome features (Fig. 2C), these phenotypes were not significantly associated on their own in the PheWAS analysis (Fig. 2D). A similar pattern was observed for the remaining variants identified in the discovery analysis (figs. S2 to S17). A PheWAS analysis of all 6188 variants tested in the discovery analysis and 1734 phecodes did not yield any significant ( $q < 0.1$ ) associations.



**Fig. 1. PheRSs capture the diagnostic fingerprint of Mendelian disease in EHR data.** Scores for six Mendelian diseases were calculated for clinically diagnosed cases and controls matched by age, sex, race, and record length. (A) Boxplots of PheRSs for cases and controls for each disease. (B) Number

of cases and statistical significance between cases and controls (Wilcoxon rank sum test) for each disease. (C) Matrix of standardized differences in location (pseudomedian) of the PheRSs between cases and controls (by row) and for each Mendelian disease definition (by column).

Replication of previously unrecognized associations

We attempted to replicate significant associations from the discovery analysis in two independent cohorts: a European ancestry cohort from Marshfield Clinic ( $n = 9441$ ) and a non-European ancestry cohort from Vanderbilt ( $n = 3820$ ) (tables S4 and S5). Each was tested as in the discovery cohort, using linear regression assuming a dominant model, adjusted for age and sex. Only variants with at least 10 heterozygotes or homozygotes for the rare allele were tested. In the Marshfield cohort, both attempted associations were replicated: p.G77S in *TG* with thyroid dysmorphogenesis PheRS ( $P = 5.0 \times 10^{-4}$ ) and p.R507H in *FAN1* with karyomegalic interstitial nephritis PheRS ( $P = 8.2 \times 10^{-3}$ ) (table S6). In the Vanderbilt non-European ancestry cohort, we replicated two of three associations: p.A993A in *KIF1A* with spastic paraplegia PheRS ( $P = 1.9 \times 10^{-3}$ ) and p.A295T in *AGXT* with primary hyperoxaluria type 1 PheRS ( $P = 3.9 \times 10^{-3}$ ). The association between p.R507H in *FAN1* and karyo-

megalic interstitial nephritis PheRS was not replicated in the non-European ancestry cohort, potentially owing to the small number of individuals with the allele in the replication cohort ( $n = 15$ ).

Sequencing individuals with hitherto unrecognized variants

To test for additional rare variants segregating with high-PheRS individuals, we analyzed the whole-exome sequences of 84 individuals from the discovery analysis for seven of the significant variants (table S7), including individuals with elevated ( $n = 36$ ) and nonelevated ( $n = 48$ ) PheRSs. A total of four individuals were found to carry a second rare, nonsynonymous variant in the target gene (Fig. 3 and tables S8 and S9). Two were possible compound heterozygotes (phase could not be determined in this analysis) (*PLCG2* and *AGXT*), and two were homozygotes for the variant identified in the discovery analysis (*DGKE* and *AGXT*), confirming the results from genotyping). Three of the four individuals with con-

firmed second variants had the highest PheRS for their respective diseases among those selected for whole-exome sequencing.

The heterozygote for *AGXT* p.A295T who was found to have an additional rare *AGXT* variant (p.R381K) through whole-exome sequencing had the highest PheRS for primary hyperoxaluria type 1, a recessive condition characterized by nephrocalcinosis and oxalate nephrolithiasis; an EHR review revealed that he had calcium oxalate crystals evident in urinalysis. The second-highest-scoring individual, a confirmed heterozygote, was diagnosed with hyperoxaluria, which was attributed to his Crohn's disease. The p.A295T homozygotes in the discovery and replication cohorts were no more symptomatic than their heterozygous counterparts. This evidence, along with the persistence of the signal after removing individuals with second variants, suggests that p.A295T may act as a strong risk factor for hyperoxaluria, with more severe manifestations occurring in individuals with additional genetic or environmental risk factors.

**Table 1. Significant associations between phenotype risk scores (PheRSs) for Mendelian disease and rare variants.** Shown are significant results from the analysis of 7520 PheRS-variant pairs, generated using linear regression assuming a dominant model, adjusted for age and sex. All associations with a false discovery rate of  $q < 0.1$  are included. The established mode of inheritance is listed in the "OMIM reported inheritance" column; AD indicates autosomal dominant, AR indicates autosomal recessive; and AR\* indicates that disease has also been reported in heterozygotes. ClinVar designations are included, when available: P, pathogenic; LP, likely pathogenic; LB, likely benign; U, uncertain significance. Variants with relevant phenotype associations in HGMD are indicated with a "Y." Results from applying American College of Medical Genetics and Genomics (ACMG) interpretations are given in the last column; an arrow indicates that the interpretation changed in light of evidence presented in this paper. rsIDs are assigned by the dbSNP database. HOM, homozygote; HET, heterozygote.

Gene	Variant	rsID	HOM/ HET	Associated Mendelian Disease	OMIM Reported inheritance	Phenotype categories in PheRS	Beta	P	ClinVar	HGMD	ACMG
<i>CFTR</i>	c.1624G>T p.Gly542Ter	rs113993959	1/27	Cystic fibrosis	AR		1.39	$2.9 \times 10^{-6}$	P	Y	P
<i>CHRNA4</i>	c.1448G>A p.Arg483Gln	rs55855125	1/21	Nocturnal frontal lobe epilepsy, 1	AD		0.58	$9.0 \times 10^{-8}$	U		U
<i>DGKE</i>	c.966G>A p.Trp322Ter	rs138924661	1/14	Nephrotic syndrome, type 7	AR		1.31	$2.8 \times 10^{-7}$	LP	Y	LP→P
<i>SUOX</i>	c.228G>T p.Arg76Ser	rs202085145	0/24	Sulfocysteinuria	AR		0.82	$1.7 \times 10^{-6}$	U		U→P
<i>CFTR</i>	c.1657C>T p.Arg553Ter	rs74597325	0/12	Cystic fibrosis	AR		1.81	$2.1 \times 10^{-6}$	P	Y	P
<i>KIF1B</i>	c.2021C>T p.Thr674Ile	rs41274468	0/21	Charcot-Marie-Tooth disease, 2A1	AD		0.79	$5.3 \times 10^{-6}$			U
<i>VWF</i>	c.5851A>G p.Thr1951Ala	rs144072210	0/21	Von Willebrand disease	AR*		0.53	$8.6 \times 10^{-6}$		Y	U
<i>KIF1A</i>	c.2676C>T p.Ala993=	rs116297894	1/25	Spastic paraplegia-30	AR		0.84	$1.3 \times 10^{-5}$	LB		LB→U
<i>F10</i>	c.872G>A p.Arg291Gln	rs149212700	0/15	Factor X deficiency	AR*		0.62	$1.9 \times 10^{-5}$			U
<i>HFE</i>	c.502G>C p.Glu168Gln	rs146519482	0/40	Hemochromatosis	AR		1.08	$4.0 \times 10^{-5}$	U	Y	U
<i>TG</i>	c.229G>A p.Gly77Ser	rs142698837	0/69	Thyroid dysmorphogenesis	AR		0.26	$6.0 \times 10^{-5}$		Y	U→P
<i>SH2B3</i>	c.1183G>A p.Glu395Lys	rs148636776	0/22	Familial erythrocytosis, 1	AD		1.48	$6.1 \times 10^{-5}$			U→P
<i>SPTBN2</i>	c.7109G>A p.Arg2370His	rs145522851	0/11	Spinocerebellar ataxia	AR*		0.75	$9.0 \times 10^{-5}$			U→LP
<i>FAN1</i>	c.1520G>A p.Arg507His	rs150393409	0/434	Interstitial nephritis, karyomegalic	AR		0.15	$9.9 \times 10^{-5}$			LB→U
<i>PANK2</i>	c.1561G>A p.Gly521Arg	rs137852959	0/26	HARP syndrome	AR		0.58	$1.1 \times 10^{-4}$	P	Y	P
<i>SH2B3</i>	c.1183G>A p.Glu395Lys	rs148636776	0/22	Essential thrombocythemia	AD		0.33	$1.4 \times 10^{-4}$			U→P
<i>AGXT</i>	c.883G>A p.Ala295Thr	rs13408961	1/35	Primary hyperoxaluria, type I	AR		0.82	$1.7 \times 10^{-4}$	U/LB		LB→U
<i>PLCG2</i>	c.751A>G p.Ile251Val	rs190840748	0/10	Familial cold autoinflammatory syn. 3	AD		0.70	$1.9 \times 10^{-4}$			U

Neoplastic

Nervous/Psychiatric/Sensory

Digestive/Genitourinary

Other symptoms/Injuries

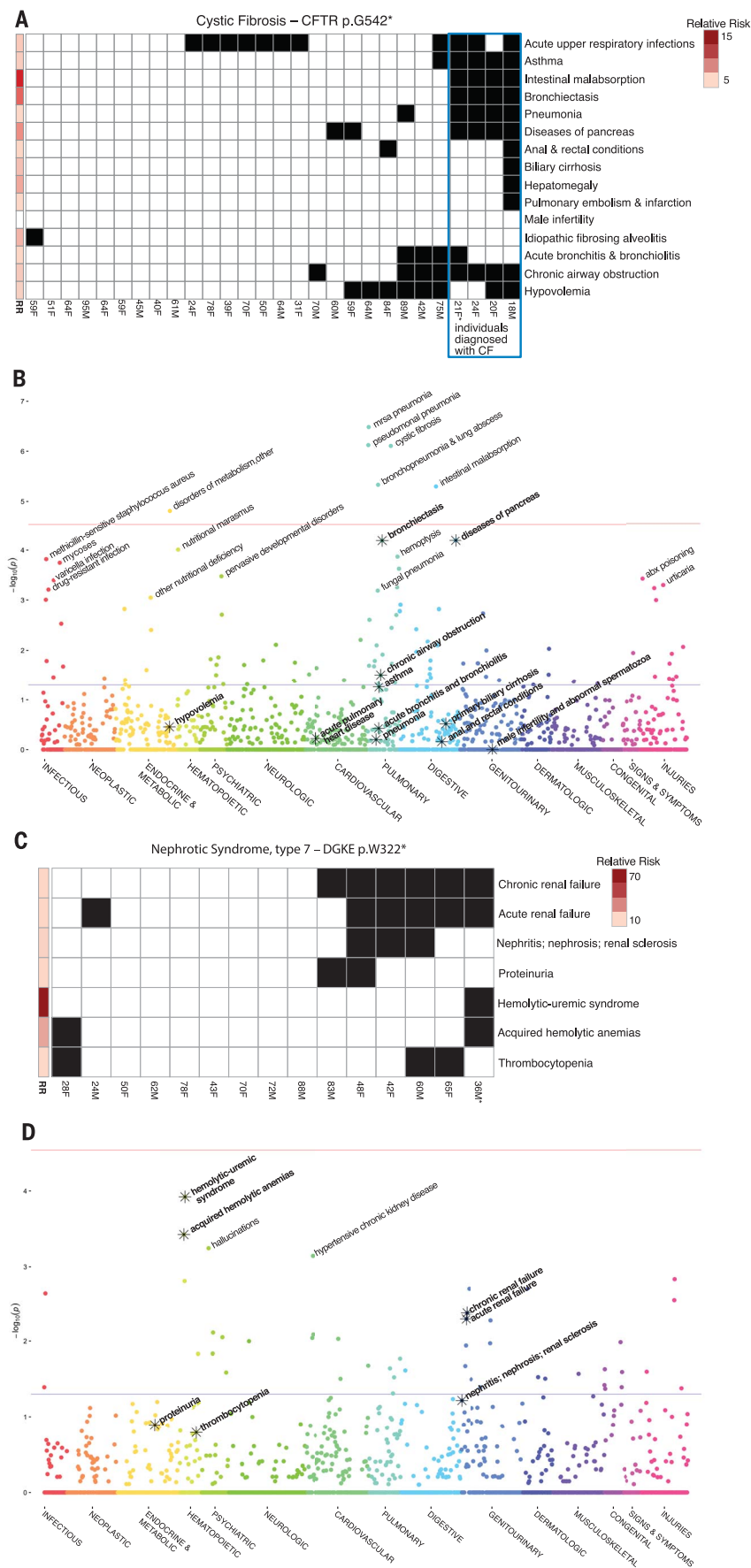
Endocrine/Metabolic/Blood

Circulatory/Respiratory

Musculoskeletal/Dermatologic

Downloaded from <http://science.sciencemag.org/> on July 30, 2018





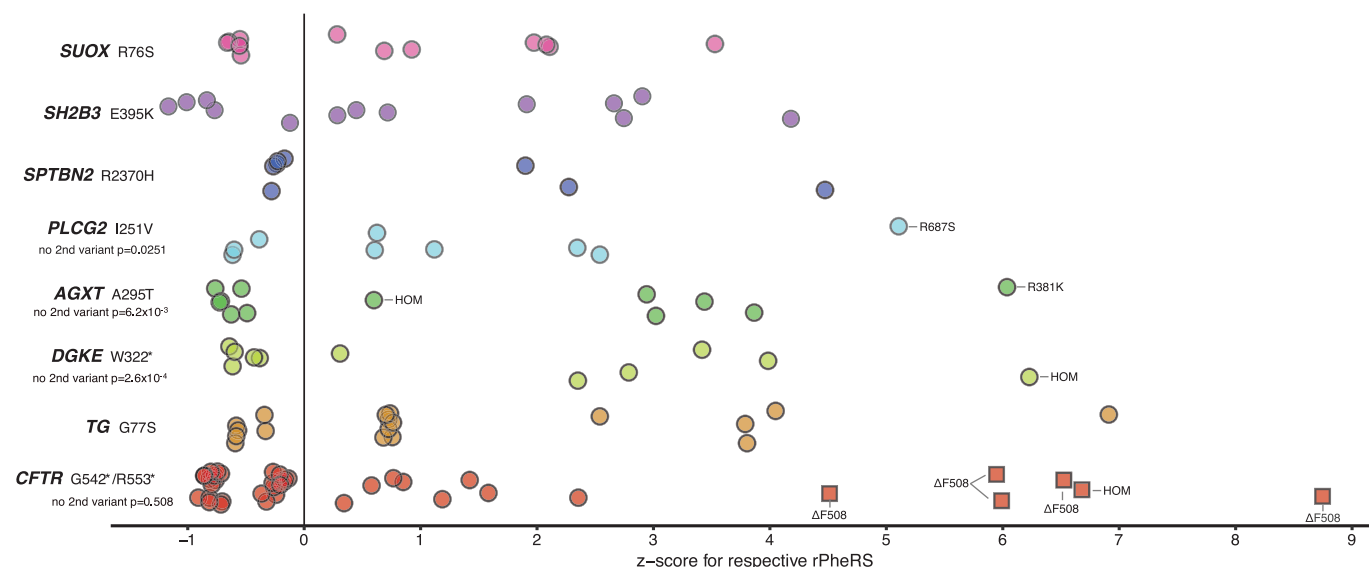
**Fig. 2. Phenotypes and PheWAS for two variants associated with PheRS for cystic fibrosis and nephrotic syndrome.** For phenotype grids (A) and (C), each row corresponds to a phenotype used in the PheRS; each column represents an individual who is heterozygous or homozygous (starred) for the variant. The bar on the left of the grid indicates the relative risk (RR) for each phenotype compared with the wild type. In grid (A), individuals clinically diagnosed with cystic fibrosis are enclosed in a blue box. PheWAS plots (B) and (D) show the PheWAS (Fisher's exact *P* value) for the variant in (A) and (C), respectively. The constituent phenotypes that define the PheRS are starred. All associations with *P* < 0.001 are labeled. The horizontal red and blue lines show the Bonferroni correction threshold for an individual PheWAS and the nominal (uncorrected) *P* = 0.05, respectively. Single-letter abbreviations for the amino acid residues are as follows: A, Ala; E, Glu; G, Gly; H, His; I, Ile; K, Lys; Q, Gln; R, Arg; S, Ser; T, Thr; V, Val; W, Trp; and Y, Tyr. An asterisk signifies a termination codon.

The highest scorer for familial cold autoinflammatory syndrome 3 (FCAS3), a dominant condition caused by variants in *PLCG2*, presented in the emergency room with a systemic “urticarial-type rash,” for which a cause was never identified, and continued to present with blistering rashes and lip and tongue swelling. Whole-exome sequencing revealed that this patient harbors a second rare variant (p.R687S) in the SH2 domain of *PLCG2*. A nearby variant (p.S707Y), also in the SH2 domain, has been implicated in a related dominant disease that has overlapping features with FCAS3 (OMIM #614878) (26).

The confirmed homozygote for p.W322\* in *DGKE*, a recessive gene that causes nephrotic syndrome type 7, was diagnosed with hemolytic uremic syndrome as a child and received a kidney transplant in his teens. A genetic etiology for his symptoms was never explored.

Sequencing confirmed the variants called on the Exome BeadChip for *SUOX*, *SH2B3*, *SPTBN2*, and *TG* and did not reveal any additional rare variants in the target genes. For *SH2B3*, the lack of a second variant is consistent with an established dominant inheritance pattern, as well as the high proportion of heterozygotes in our cohort (20 of 22) with at least one feature of familial erythrocytosis. Individuals without a second variant in *AGXT* and *DGKE*, both of which are associated with recessive conditions, also had elevated PheRSs (Fig. 3). This stands in contrast to the heterozygotes for the *CFTR* variants, who did not have a significantly elevated PheRS (*P* = 0.51; linear regression assuming a dominant model, adjusted for age and sex), consistent with a recessive inheritance model. These findings suggest a blurring of the distinction between dominant and recessive labels for some genes.

Sequencing did not reveal any additional rare nonsynonymous variants in the 36 individuals with nonelevated PheRSs. In general, individuals with the highest PheRSs were more likely to be



**Fig. 3. Whole-exome sequencing reveals second variants among individuals with high PheRSs and demonstrates disease risk in heterozygotes.** Each point represents an individual who is heterozygous or homozygous for the variant labeled on the left. The x axis represents the z-score for the PheRS relative to what is expected given age and sex (using the residual from the PheRS). All individuals carry at least one copy of the variant indicated on the left; additional variants identified by whole-exome sequencing or clinical chart review are labeled for each individual;

homozygotes confirmed by sequencing are labeled “HOM.” Additional *CFTR* variants were ascertained from clinical testing in the EHR; all other individuals were sequenced for this study. Clinically diagnosed individuals are represented as squares; all others are shown as circles. Where additional variants were found, the association test from the discovery analysis was repeated after dropping individuals with a second variant (*P* values generated using linear regression assuming a dominant model, adjusted for age and sex), and the *P* value is recorded under the gene/variant label.

clinically diagnosed or have additional genetic variants related to their symptoms (fig. S18).

### Biologic validation of *SH2B3*, *TG*, and *SUOX* associations

We selected three candidate previously unrecognized associations for biologic validation: *SH2B3*, *SUOX*, and *TG*. *SH2B3* is a negative regulator of cytokine signaling in hematopoietic cells that operates through a direct interaction between its SH2 domain and JAK2 to attenuate JAK2-mediated activation of proliferative pathways (27). The variant identified in this study, p.E395K, is located in a region of the protein that is critical for its inhibitory function (28) and is near known disruptive variants (29). Human embryonic kidney (HEK) 293T cells stimulated with erythropoietin showed an increase in phosphorylated extracellular signal-regulated kinase (pERK) levels that was quenched in the presence of wild-type *SH2B3* but not quenched with both the known p.R392E functional mutation and our p.E395K variant (Fig. 4, A and B).

Splicing prediction programs suggested a probable reduction in 5' donor strength for *SUOX* p.R76S and possible generation of an exonic cryptic splice acceptor site by *TG* p.G77S. *SUOX* p.R76S is located at the conserved -1 position of the 5' donor of exon 5. We demonstrated that the *SUOX* variant caused a decrease in exon inclusion from 96 to 35% (unpaired two-tailed *t* test; *P* < 0.001; Fig. 4C). No transcripts aside from the exon-included and exon-skipped transcripts were detected. Similarly, *TG* p.G77S resulted in altered splicing. The basal rate of exon inclusion was re-

duced from 65% for the wild-type *TG* exon to only 26% inclusion in the p.G77S exon (unpaired two-tailed *t* test; *P* < 0.001). These ratios were consistent across a range of cDNA concentrations and polymerase chain reaction (PCR) cycle numbers (Fig. 4, C and D).

### Comparison of PheRS with existing methods to determine variant pathogenicity

Across all PheRS variant associations with nominal *P* < 0.05 (*n* = 454), functional annotations were significantly correlated with PheRS effect size (Wilcoxon rank sum test); splice donor/acceptor and stop-gain variants tended to have the largest effect size, followed in decreasing order by missense, splice region, synonymous, and intron/untranslated region variants (fig. S19A). Thirteen of 14 functional prediction methods trended or associated with the probability of finding associations with the PheRS; predictions from CADD (combined annotation-dependent depletion) (9), SiPhy (30), and Polyphen2 HVAR (10) were statistically significant (*P* < 0.05; Fisher's exact test; fig. S19B).

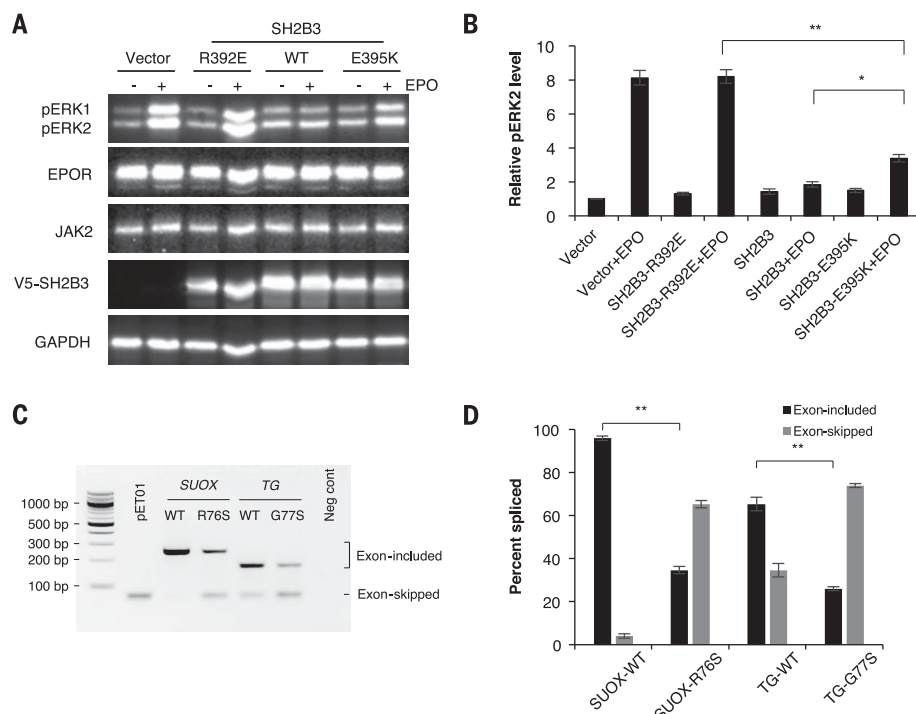
### Discussion

In our validation study, PheRS was very effective in identifying patients with diagnosed Mendelian disease by using only the phenotypic signatures. Applying PheRS to a genotyped population, we found an increased burden of phenotypes among individuals with rare variants in Mendelian disease genes. Sequencing identified or confirmed second rare variants in four individuals, three of

whom had the highest PheRS among all heterozygotes or homozygotes for that variant. In vitro studies provided supporting evidence of pathogenicity for all three variants tested.

Although our approach relies on many decades of accumulated knowledge about the phenotypic imprint of Mendelian disease, the method itself is simple to implement. Our ability to replicate results in an external cohort suggests that it is portable and would therefore be applicable to data sets such as those of the Million Veteran Program, UK Biobank, and the All of Us Research Program (All of Us is a service mark of the U.S. Department of Health and Human Services). Applied to such large populations, this method could facilitate the discovery of pathogenic variants, refine estimates of penetrance across diverse populations, and provide a more nuanced understanding of inheritance patterns, which this study suggests may be more complex than merely “recessive” or “dominant” for some genes. Incorporation of richer EHR data, such as laboratory results and clinical notes (31), could increase the resolving power of PheRS. Furthermore, this method may be used with other combinations of phenotypes that do not follow established Mendelian patterns, perhaps based on undiagnosed patients with unusual presentations.

The American College of Medical Genetics and Genomics established guidelines for variant interpretation that reflect the need to combine multiple lines of evidence, including population-based genotype-phenotype correlations (32). Our method provides a high-throughput means to generate such evidence. Using these guidelines, 10 of the



**Fig. 4. PheRS enriches for variants with altered function in vitro.** Representative Western blots (A) and mean pERK2 levels normalized to erythropoietin receptor (EPOR) expression (B) in EPO-stimulated HEK293T cells transiently transfected with wild-type (WT) or variant SH2B3 constructs, EPOR, and JAK2. As expected, known functional mutation SH2B3-R392E fails to inhibit EPO-stimulated ERK phosphorylation. Similarly, SH2B3-E395K shows ~1.8-fold elevation (relative to WT SH2B3) of EPO-stimulated ERK activation at 10 min. Reverse transcription PCR analysis (C) and quantification (D) of WT versus variant splicing of *SUOX* and *TG* exons in HEK293T cells transiently transfected with empty minigene vector pET01, pET01 containing exons of interest flanked by 100 base pairs (bp) of intronic sequence, or negative control pIRES2-EGFP (enhanced green fluorescent protein). Absolute change in the percent of exon inclusion was ~61% for the *SUOX* variant and ~39% for the *TG* variant. Means  $\pm$  SEM;  $n = 4$  [(A) and (B)] or 5 [(C) and (D)] independent experiments; unpaired two-tailed  $t$  test; \* $P = 0.003$ , \*\* $P < 0.001$ .

variants from the discovery analysis were interpreted as having “uncertain significance.” By adding data from the PheRS analysis, in combination with evidence from our in vitro studies, four of these variants could be converted to “likely pathogenic” or “pathogenic” (Table 1 and table S10).

Our findings suggest that the phenotypic burden of rare variants in Mendelian genes may be greater than previously thought. A combination of PheRS and sequencing identified symptomatic individuals with genetics consistent with established inheritance patterns—heterozygous individuals for dominant genes (*SH2B3* and *PLCG2*) and individuals with confirmed second variants in recessive genes (*DGKE* and *AGXT*)—none of whom were diagnosed with a genetic condition. A much larger number of individuals were heterozygous for variants in genes with a presumed recessive inheritance and yet still had symptoms consistent with the Mendelian disease pattern. Although we cannot exclude the possible influence of structural or noncoding variants, the evidence suggests that these variants increase risk in heterozygotes. These individuals tend to have disease that is mild compared with the classic

presentations but severe relative to the general population. For example, homozygous pathogenic mutations in *TG* are associated with congenital goiter, which often progresses to thyroid carcinoma; our most severely affected heterozygote received a thyroidectomy at age 26 for goiter and thyroid carcinoma.

This work adds to the evidence that Mendelian and complex disease are not dichotomous, but rather exist on a spectrum. As a method that is both high-throughput and sensitive to the vast knowledge already acquired, PheRS is a tool that may help bridge the gap between Mendelian and complex disease. A consequential question is whether the treatments designed for a Mendelian condition could be effective in individuals with nontraditional molecular presentations. Of the 17 diseases represented among those patients with suspected but undiagnosed Mendelian disease, 11 have specific treatments available (table S11), some of which could alter the long-term course of the disease.

The impact that this approach will have on accelerating precision medicine depends on three interrelated challenges. First, we must integrate

statistical associations generated with PheRS into guidelines used for variant interpretation. Second, as we collect stronger evidence for the phenotypic effects of rare variants, we must learn to rapidly and effectively integrate that knowledge into clinical care. Third, we must determine whether PheRS can be used to prospectively identify patients whose symptoms are caused by variants in Mendelian genes. If these challenges are addressed, approaches like ours may ultimately enable the conversion of big data not just to knowledge but also to improved care and outcomes for patients.

## REFERENCES AND NOTES

- J. R. Lupski, J. W. Belmont, E. Boerwinkle, R. A. Gibbs, *Cell* **147**, 32–43 (2011).
- D. R. Blair et al., *Cell* **155**, 70–80 (2013).
- T. Groza et al., *Am. J. Hum. Genet.* **97**, 111–124 (2015).
- D. Langlais, N. Fodil, P. Gros, *Annu. Rev. Immunol.* **35**, 1–30 (2017).
- J. C. Denny, L. Bastarache, D. M. Roden, *Annu. Rev. Genomics Hum. Genet.* **17**, 353–373 (2016).
- W. S. Bush, M. T. Oetjens, D. C. Crawford, *Nat. Rev. Genet.* **17**, 129–145 (2016).
- E. Marouli et al., *Nature* **542**, 186–190 (2017).
- D. G. MacArthur et al., *Nature* **508**, 469–476 (2014).
- M. Kircher et al., *Nat. Genet.* **46**, 310–315 (2014).
- I. A. Adzhubei et al., *Nat. Methods* **7**, 248–249 (2010).
- N. M. Ioannidis et al., *Am. J. Hum. Genet.* **99**, 877–885 (2016).
- H. L. Rehm et al., *N. Engl. J. Med.* **372**, 2235–2242 (2015).
- M. Lek et al., *Nature* **536**, 285–291 (2016).
- A. K. Manrai et al., *N. Engl. J. Med.* **375**, 655–665 (2016).
- R. Chen et al., *Nat. Biotechnol.* **34**, 531–538 (2016).
- S. L. Van Driest et al., *JAMA* **315**, 47–57 (2016).
- I. S. Kohane, *Nat. Rev. Genet.* **12**, 417–428 (2011).
- D. C. Crawford et al., *Front. Genet.* **5**, 184 (2014).
- OMIM (Online Mendelian Inheritance in Man), <http://omim.org/>.
- S. Köhler et al., *Nucleic Acids Res.* **42**, D966–D974 (2014).
- J. C. Denny et al., *Nat. Biotechnol.* **31**, 1102–1110 (2013).
- A. Verma et al., *PLOS ONE* **11**, e0160573 (2016).
- W.-Q. Wei et al., *PLOS ONE* **12**, e0175508 (2017).
- E. L. Macleod, D. M. Ney, *Ann. Nestlé [Engl.]* **68**, 58–69 (2010).
- P. D. Stenson et al., *Hum. Genet.* **136**, 665–677 (2017).
- Q. Zhou et al., *Am. J. Hum. Genet.* **91**, 713–720 (2012).
- N. Maslah, B. Cassinat, E. Verger, J.-J. Kiladjian, L. Velazquez, *Leukemia* **31**, 1661–1670 (2017).
- W. Tong, J. Zhang, H. F. Lodish, *Blood* **105**, 4604–4612 (2005).
- C. Camps et al., *Haematologica* **101**, 1306–1318 (2016).
- M. Garber et al., *Bioinformatics* **25**, i54–i62 (2009).
- S. J. Hebbaring et al., *Bioinformatics* **31**, 1981–1987 (2015).
- S. Richards et al., *Genet. Med.* **17**, 405–423 (2015).

## ACKNOWLEDGMENTS

The authors thank B. Li and Q. Wei for technical assistance with the whole-exome sequence variant-calling pipeline. **Funding:** This work was supported by grants R01-LM010685, K22-LM011939, and T15-LM007359 from the National Library of Medicine; U01 grants supporting Vanderbilt’s participation in the eMERGE (Electronic Medical Records and Genomics) network (HG004603, HG006378, and HG008672); and grants T32-HG008341 and U01-HG009086 from National Human Genome Research Institute. BioVU received and continues to receive support through the National Center for Research Resources (UL1-RR024975), which is now the National Center for Advancing Translational Sciences (UL1-TR000445). Other support comes from grants R01-GM114128, P50-GM115305, and R01-GM120523 from the National Institute for General Medical Sciences; R01-HL133786 from the National Heart, Lung, and Blood Institute; an American Heart Association career development award (I6FTF30130005); and research funds from the University of Oklahoma Health Sciences Center. **Author contributions:** L.B., J.J.H., and J.C.D. conceived of and implemented PheRS. M.B. and S.H. assisted with the replication of primary results. L.B., R.H., T.L.M., J.C., and J.C.D. researched and interpreted rare variants. J.M., J.C., and A.G. designed and conducted the *SUOX* and *TG* in vitro studies. W.Z., W.T.H., and Z.J.Z. designed and conducted the SH2B3 in vitro

studies. W.-Q.W., Q.F., D.R., T.E., and D.R.V.E. assisted with designing and analyzing the sequence data. T.E. and J.D.M. provided support for statistical analyses. T.O., A.H.R., J.M., L.B., and J.C.D. reviewed medical records. S.L.V.D., J.D.M., Q.S.W., M.T., R.C., N.J.C., and D.M.R. advised on the interpretation of the results and methods refinement. L.B., J.J.H., D.M.R., and J.C.D. drafted the manuscript. All authors provided critical feedback and helped shape the analysis and interpretation of the results of this study. **Competing interests:** S.L.V.D. received an honorarium from Merck as an invited speaker. After completing

her work on this project, T.L.M. began working for Alnylam Pharmaceuticals and received stock in the company. Neither Merck nor Alnylam supported or were involved in this research. The remaining coauthors have no competing interests to report. **Data and materials availability:** Software for PheWAS analysis and visualization is available at <https://github.com/PheWAS/PheWAS>. Software to generate PheRSs and produce grid visualizations is available at <https://github.com/labastar/PheRS>. Genetic and phenotypic data used in this study are available at dbGap under accession number phs001516.v1.p1.

#### SUPPLEMENTARY MATERIALS

[www.sciencemag.org/content/359/6381/1233/suppl/DC1](http://www.sciencemag.org/content/359/6381/1233/suppl/DC1)  
Materials and Methods  
Figs. S1 to S20  
Tables S1 to S17  
References (33–44)

14 November 2016; resubmitted 25 August 2017  
Accepted 22 January 2018  
10.1126/science.aal4043

## Phenotype risk scores identify patients with unrecognized Mendelian disease patterns

Lisa Bastarache, Jacob J. Hughey, Scott Hebring, Joy Marlo, Wanke Zhao, Wanting T. Ho, Sara L. Van Driest, Tracy L. McGregor, Jonathan D. Mosley, Quinn S. Wells, Michael Temple, Andrea H. Ramirez, Robert Carroll, Travis Osterman, Todd Edwards, Douglas Ruderfer, Digna R. Velez Edwards, Rizwan Hamid, Joy Cogan, Andrew Glazer, Wei-Qi Wei, QiPing Feng, Murray Brilliant, Zhizhuang J. Zhao, Nancy J. Cox, Dan M. Roden and Joshua C. Denny

*Science* **359** (6381), 1233-1239.  
DOI: 10.1126/science.aal4043

### Hidden effects of Mendelian inheritance

Identifying the determinate factors of genetic disease has been quite successful for Mendelian inheritance of large-effect pathogenic variants. In these cases, two non- or low-functioning genes contribute to disease. However, Mendelian effects of lesser strength have generally been ignored when looking at genomic consequences in human health. Bastarache *et al.* used electronic records to identify the phenotypic effects of previously unidentified Mendelian variations. Their analysis suggests that individuals with undiagnosed Mendelian diseases may be more prevalent in the general population than assumed. Because of this, genetic analysis may be able to assist clinicians in arriving at a diagnosis.

*Science*, this issue p. 1233

#### ARTICLE TOOLS

<http://science.sciencemag.org/content/359/6381/1233>

#### SUPPLEMENTARY MATERIALS

<http://science.sciencemag.org/content/suppl/2018/03/14/359.6381.1233.DC1>

#### REFERENCES

This article cites 43 articles, 5 of which you can access for free  
<http://science.sciencemag.org/content/359/6381/1233#BIBL>

#### PERMISSIONS

<http://www.sciencemag.org/help/reprints-and-permissions>

Use of this article is subject to the [Terms of Service](#)