



## ARTICLE

# Commonalities across computational workflows for uncovering explanatory variants in undiagnosed cases

Shilpa Nadimpalli Kobren<sup>1</sup>, Dustin Baldrige<sup>2</sup>, Matt Velinder<sup>3</sup>, Joel B. Krier<sup>4</sup>, Kimberly LeBlanc<sup>1</sup>, Cecilia Esteves<sup>1</sup>, Barbara N. Pusey<sup>5</sup>, Stephan Züchner<sup>6</sup>, Elizabeth Blue<sup>7</sup>, Hane Lee<sup>8,9</sup>, Alden Huang<sup>8</sup>, Lisa Bastarache<sup>10</sup>, Anna Bican<sup>10</sup>, Joy Cogan<sup>10</sup>, Shruti Marwaha<sup>11</sup>, Anna Alkelai<sup>12</sup>, David R. Murdock<sup>13</sup>, Pengfei Liu<sup>13,14</sup>, Daniel J. Wegner<sup>2</sup>, Alexander J. Paul<sup>15</sup> and Undiagnosed Diseases Network\*, Shamil R. Sunyaev<sup>1,4</sup> and Isaac S. Kohane<sup>1</sup>✉

**PURPOSE:** Genomic sequencing has become an increasingly powerful and relevant tool to be leveraged for the discovery of genetic aberrations underlying rare, Mendelian conditions. Although the computational tools incorporated into diagnostic workflows for this task are continually evolving and improving, we nevertheless sought to investigate commonalities across sequencing processing workflows to reveal consensus and standard practice tools and highlight exploratory analyses where technical and theoretical method improvements would be most impactful.

**METHODS:** We collected details regarding the computational approaches used by a genetic testing laboratory and 11 clinical research sites in the United States participating in the Undiagnosed Diseases Network via meetings with bioinformaticians, online survey forms, and analyses of internal protocols.

**RESULTS:** We found that tools for processing genomic sequencing data can be grouped into four distinct categories. Whereas well-established practices exist for initial variant calling and quality control steps, there is substantial divergence across sites in later stages for variant prioritization and multimodal data integration, demonstrating a diversity of approaches for solving the most mysterious undiagnosed cases.

**CONCLUSION:** The largest differences across diagnostic workflows suggest that advances in structural variant detection, noncoding variant interpretation, and integration of additional biomedical data may be especially promising for solving chronically undiagnosed cases.

*Genetics in Medicine* (2021) 23:1075–1085; <https://doi.org/10.1038/s41436-020-01084-8>

## INTRODUCTION

Next-generation exome sequencing (ES) and genome sequencing (GS) have revolutionized the process for diagnosing rare and novel genetic conditions.<sup>1</sup> Traditionally, the diagnostic process was primarily driven by phenotype, with clinicians comparing patients' symptoms to others encountered in their prior experience and clinical training and/or to a knowledgebase of known human diseases.<sup>2</sup> In a typical undiagnosed case, however, either a patient's phenotype is not indicative of any known disease, or tests to confirm the presence of a suspected genetic condition are inconclusive. In these instances, ES and GS have enabled health-care providers to pursue a genetics-driven diagnostic approach in parallel, where the genetic variation uncovered in a patient can be assessed with respect to not only its known phenotypic associations<sup>3</sup> but also to its prevalence in background populations,<sup>4</sup> predicted pathogenicity,<sup>5</sup> functional consequences, and mode of inheritance to reveal novel disease-causing loci. Indeed, while traditional clinical case review and directed diagnostic assays continue to solve difficult cases, ~74% of newly diagnosed genetic conditions have been attributed to analyses of ES and GS data.<sup>6,7</sup>

However, the diagnosis rate for patients with potentially unique genetic conditions is still ~35%,<sup>7</sup> suggesting ample opportunity for methodological improvements to advance our understanding of the genetic underpinnings of phenotypic extremes.

With this goal in mind, cross-institutional initiatives such as Care4Rare in Canada (<http://care4rare.ca>) and Solve-RD in Europe (<http://solve-rd.eu>) have been established to connect and enable clinical researchers to uncover the genetic origins of disease in undiagnosed patients. In addition to furthering basic genetics research, these efforts have provided scores of patients with an end to diagnostic uncertainty and access to additional services.<sup>8</sup> The most expansive undiagnosed initiative in the United States is the Undiagnosed Diseases Network (UDN), which encompasses 12 clinical sites and has, since its inception in 2014, cumulatively diagnosed over 400 individuals and described over 30 novel syndromes.<sup>7</sup> Each UDN clinical site is staffed with specialists who develop and apply complex suites of bioinformatics tools to analyze sequencing data and uncover disease-causing variants.<sup>9</sup> These sites each underwent a competitive application process and were selected to join the UDN due to their demonstrated track

<sup>1</sup>Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA. <sup>2</sup>Department of Pediatrics, Washington University School of Medicine, St. Louis, MO, USA.

<sup>3</sup>Center for Genomic Discovery, University of Utah, Salt Lake City, UT, USA. <sup>4</sup>Division of Genetics, Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA.

<sup>5</sup>National Human Genome Research Institute (NHGRI) at the National Institutes of Health (NIH), Bethesda, MD, USA. <sup>6</sup>Department of Human Genetics and Hussman Institute for Human Genomics, University of Miami Health System, Miami, FL, USA. <sup>7</sup>Division of Medical Genetics, Department of Medicine, University of Washington, Seattle, WA, USA.

<sup>8</sup>Department of Human Genetics, David Geffen School of Medicine at the University of California, Los Angeles, CA, USA. <sup>9</sup>Department of Pathology and Laboratory Medicine, David Geffen School of Medicine at the University of California, Los Angeles, CA, USA. <sup>10</sup>Department of Biomedical Informatics, Vanderbilt University Medical Center, Nashville, TN, USA. <sup>11</sup>Stanford Center for Undiagnosed Diseases, Stanford, CA, USA. <sup>12</sup>Institute for Genomic Medicine, Columbia University Medical Center, New York City, NY, USA.

<sup>13</sup>Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, TX, USA. <sup>14</sup>Baylor Genetics, Houston, TX, USA. <sup>15</sup>McDonnell Genome Institute, Washington University School of Medicine, St. Louis, MO, USA. \*A list of authors and their affiliations appears at the end of the paper. ✉email: isaac\_kohane@hms.harvard.edu

record of diagnosing difficult cases and characterizing novel genetic conditions through ongoing research efforts. The workflows implemented at these sites are thus representative of the state-of-the-art in rare disease diagnostic efforts.

We gathered details about 12 UDN bioinformatics pipelines, determined recurrent steps in a typical diagnostic evaluation, and identified consensus approaches. Moreover, we highlight substantial differences across pipelines regarding overall organization and incorporated tools. The comprehensive snapshot of effective computational workflows presented here can direct clinical teams interested in initiating genomic sequencing usage or re-evaluating patients who have had inconclusive genetic testing.

## MATERIALS AND METHODS

### Participating sites

Sequence analysis pipeline details were collected from the CLIA-certified sequencing core at Baylor Genetics (BaylorSeq) and 11 UDN clinical sites: Baylor College of Medicine (BCM), Duke University and Columbia University Institute for Genomic Medicine (Duke/Columbia), three Harvard-affiliated hospitals and Brigham Genomic Medicine (Harvard), University of Miami Miller School of Medicine (Miami), National Institutes of Health (NIH), University of Washington School of Medicine and Seattle Children's Hospital (PacificNW), Stanford Center for Undiagnosed Diseases (Stanford), University of California–Los Angeles (UCLA), University of Utah Health Center for Genetic Discovery (Utah), Vanderbilt University Medical Center (Vanderbilt), and Washington University School of Medicine (WUSTL). The University of Pennsylvania and Children's Hospital of Philadelphia clinical site had yet to process sequencing data for a UDN case at the time of writing and thus is excluded from this study.

### Data collection

We systematically collected details about each UDN site's computational diagnostic workflows using a combination of in-person and virtual meetings with bioinformaticians and genetic counselors, online survey forms, and inspections of published papers and internal protocols.<sup>10–12</sup>

## RESULTS

### Overview of diagnostic workflow components

Before applying to the UDN, a patient has typically endured extensive prior testing by multiple clinicians over the course of a multiyear “diagnostic odyssey.” As part of the application process, UDN clinical sites review patients' health records to assess whether the UDN evaluation may aid in the identification of a diagnosis. Accepted patients undergo an in-person evaluation at a clinical site (Fig. 1a). In most cases, blood, saliva, and/or fibroblast samples of affected and unaffected individuals in the family are collected during this evaluation or beforehand via mailed-in collection kits. These samples are sequenced at BaylorSeq; all sequencing data are made available to the clinical site within weeks (Fig. 1b). Variants in disease-causing genes related to the clinical phenotype, medically actionable pathogenic variants in disease-causing genes unrelated to the clinical phenotype, and heterozygote status for select recessive Mendelian conditions are listed in a clinical report issued by BaylorSeq in accordance with the UDN protocol and following American College of Medical Genetics and Genomics (ACMG) variant classification guidelines.<sup>13</sup> At 8 of the 11 clinical sites surveyed, researchers simultaneously perform local analyses of the sequencing data in an attempt to identify “strong candidate” variants that may explain the patient's symptoms (Fig. 1c, d); three surveyed sites run their local pipelines only when BaylorSeq's clinical report is inconclusive. Once candidate variants are highlighted via clinical sites' and BaylorSeq's analyses, there are three ways by which their causality is established. First, human and animal databases are queried for genotype-matched individuals with symptomatic concordance with the patient.<sup>14–17</sup> Second, experiments are simultaneously performed to evaluate the in vivo effect of candidate

variants in model organisms or cell lines. Third, the presence of secondary phenotypes indicated by genotype-matched individuals or in vivo experiments are confirmed in affected patients (Fig. 1e). Causal variants revealed through these steps are confirmed by Sanger sequencing, broadly shared by the UDN (Extended Data Note 1), and ideally lead to a molecular diagnosis for a patient, which in and of itself represents a turning point in a patient's diagnostic odyssey, and also can inform positive therapeutic changes (Fig. 1f).<sup>18</sup>

The computational tools used to find explanatory genetic variants change constantly with newly available technologies and newly encountered disease etiologies. Despite these iterative improvements to bioinformatics pipelines, the primary roles that computational tools play in the overall variant prioritization process can be categorized as follows: (1) aligning sequencing reads to a reference human genome (Fig. 1g), (2) identifying genetic variants present in the individual from the sequencing reads (Fig. 1h), (3) annotating those variants with relevant information (Fig. 1i), and finally (4) filtering and prioritizing variants that are likely to cause the patient's condition (Fig. 1j). In the following sections, we delve into the purpose of and tools used in each of these categories.

### Aligning next-generation sequencing reads

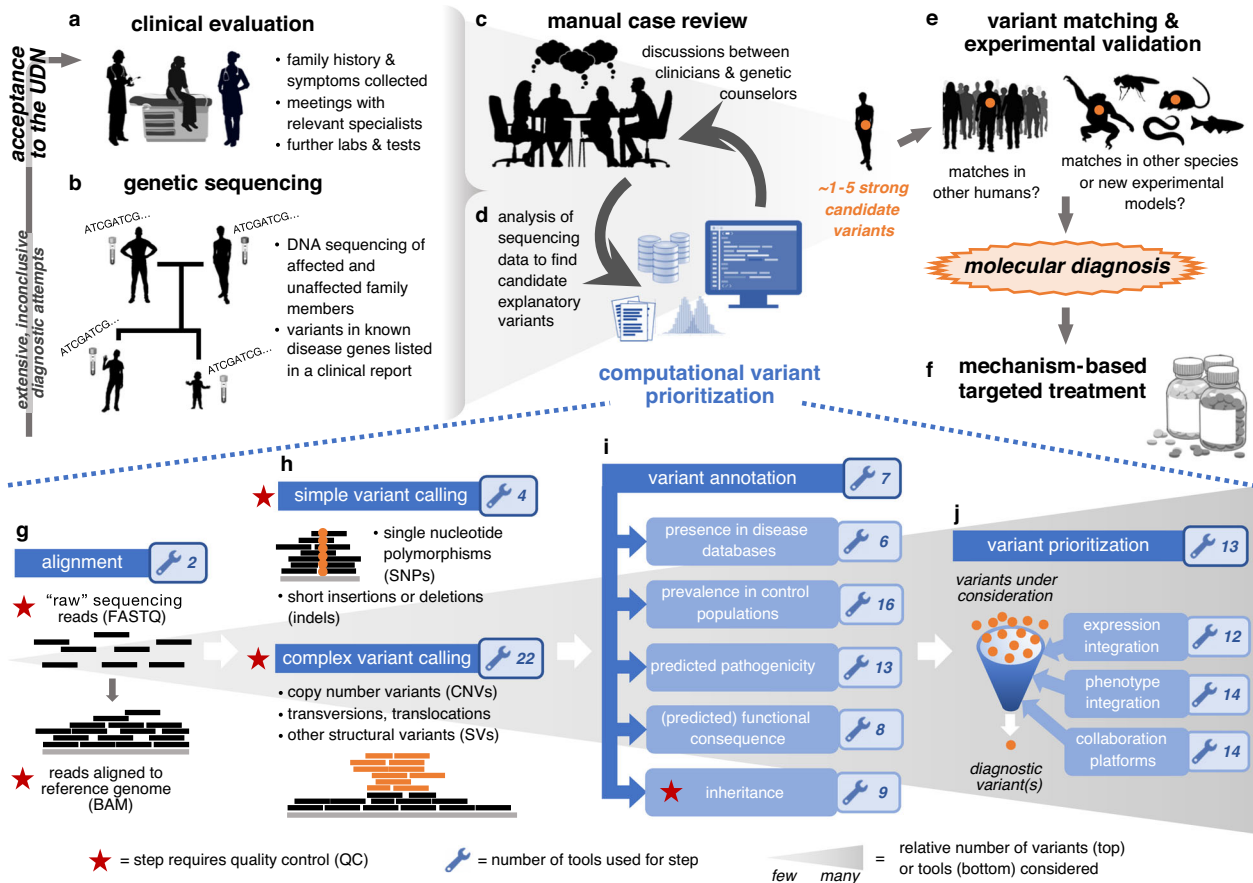
Aligning next-generation sequencing reads to a reference human genome is the necessary first step for all sequence analysis pipelines (Fig. 1g); the ubiquity of this step has resulted in community-driven standardization.<sup>19</sup> Eight sites regularly realign reads after BaylorSeq's initial alignment, whereas three sites realign reads only in specific circumstances, such as during reanalysis of a patient's prior sequencing data. Realignment is necessary for six sites whose pipelines are configured for the GRCh37/hg19 human genome build, as genetic testing laboratories including BaylorSeq now provide reads aligned to the newer GRCh38/hg38 build. Realignment uses either an open-source implementation of the Burrows–Wheeler Aligner (BWA-MEM) (used regularly by six sites and in specific circumstances, as described above, by two sites) or Illumina/Edico's DRAGEN aligner (used regularly by BaylorSeq and two clinical sites and in specific circumstances by one clinical site).

### Simple variant calling

Calling single-nucleotide variants (SNVs) and short insertions and deletions (indels) from aligned reads is the next step in sequence processing (Fig. 1h) and is often accomplished using the Genome Analysis Toolkit (GATK) best practices workflow,<sup>20</sup> though Google's DeepVariant<sup>21</sup> and Real Time Genomics' PolyBayes implementation (<https://www.realtimegenomics.com>) perform competitively for this task and are used in addition to GATK by two clinical sites. BaylorSeq calls variants using Illumina/Edico's DRAGEN platform. Six clinical sites and BaylorSeq “jointly” call variants across samples as recommended in GATK to rescue low coverage true variants and accurately model false variants. In practice, variants are jointly called with (1) members of the same family, (2) other UDN patients at the same site, and/or (3) healthy patients internal or external to an institution. The Variant Quality Score Recalibration (VQSR) step recommended by GATK to identify technical artifacts, however, may misclassify real rare variants as false positives; this step is carefully reviewed or omitted in practice.

### Structural variant detection

In contrast to calling simple variants, calling structural variants (SVs) from GS data is a relatively divergent step, indicating that best practices have yet to be determined. SVs refer to large (>50 bp) insertions and deletions, duplications and other copy-number variants (CNVs), short tandem repeat (STR) expansions, translocations where genomic regions have moved within or



**Fig. 1 Representative clinical workflow to uncover disease-causing genetic variants in undiagnosed patients.** Upon acceptance to the Undiagnosed Diseases Network (UDN), (a) an affected patient has an in-person clinical evaluation where extensive phenotyping and additional tests are performed as needed. (b) Before or during the clinical evaluation, samples of relevant affected and unaffected individuals in a family are sent for genomic sequencing. (c,d) Sequencing data provided by the sequencing center are analyzed in conjunction with other information in a back-and-forth process between bioinformaticians, clinicians, and genetic counselors to highlight variants that are likely to explain the patient's disease. (e) Matches to the strong candidate explanatory variants identified in (c) are searched for in databases containing human genetic variant and corresponding symptom information (e.g., Matchmaker Exchange) or in databases containing animal genetic variants and corresponding phenotype information (e.g., MARRVEL). Strong candidate variants are also introduced into model organisms or cell lines where possible to assess in vivo phenotypic impact. (f) Once a candidate variant has been confirmed as disease causal, a molecular diagnosis is provided that can subsequently be used to tailor clinical management and molecular therapeutics. (g–j) Recurring steps in computational workflows to process genomic sequencing data to call, filter, and prioritize genetic variants that explain the affected individual's disease symptoms.

across chromosomes, and inversions where a detached stretch of DNA was reattached in the opposite orientation. Combining the output from many SV calling tools—each optimized for detecting complementary types of SVs and often using distinct information (e.g., read depth, paired-end reads, or split reads)—is necessary for comprehensive SV detection.<sup>22</sup> Existing SV detection tools have been reviewed in depth;<sup>23</sup> here we list the subset of tools that are actively used by UDN sites (Table 1, Extended Data Table 1). The most commonly used tool, Manta, has been shown by independent evaluations to have high sensitivity but also a high false positive rate.<sup>24</sup> Future development of SV benchmarking data sets for assessing the accuracy of SV detection tools will be essential in directing the current diverse exploration of techniques toward community-established best practices.

#### Quality control of called variants

Confirming the quality of sequencing data and variants is critical to avoid expending downstream analyses on false variants. CLIA-certified genetic testing laboratories check the quality of unaligned and map-aligned sequencing reads prior to variant calling for all clinical grade sequencing (Extended Data Note 2). Four UDN clinical sites regularly confirm the quality of sequencing

reads using a combination of FASTQC, FASTP, MultiQC, BEDTools (to check coverage), and bam.iobio. Other clinical sites begin quality control (QC) only after read alignment and variant calling.

QC for Mendelian disease diagnosis encompasses three checks: (1) sequencing reads are high quality, (2) sequenced samples correspond to the correct individuals and have expected relatedness, and (3) inheritance patterns across families are as expected (Table 2, Extended Data Table 1). BaylorSeq performs QC for all clinical genomic sequencing before providing data to UDN clinical sites. However, when patients provide their own sequencing data (as opposed to BaylorSeq providing newly acquired data) or when "research" (as opposed to clinical) sequencing is provided, clinical sites perform QC. Most sites have nearly identical steps for check 1 and similar QC for checks 2 and 3. In practice, QC has identified incorrectly related or labeled samples and poor overall quality of sequencing reads that were remedied via resequencing before subsequent analyses.<sup>11</sup> Notably, existing QC tools rarely "flag" anomalous samples; users must accurately interpret results.

#### Annotation and filtering of genetic variants

Even after removing low quality calls, a single genome can have several thousand unique genetic variants uncovered. Efficient,

**Table 1.** Structural variant (SV) callers in use at clinical sites.

	BaylorSeq	BCM	Duke/ Columbia	Harvard	Miami	NIH	PacificNW	Stanford	UCLA	Utah	Vanderbilt	WUSTL
<b>Find SVs from sequencing reads</b>												
Manta <sup>a</sup>	■	■	□	□	□	□	□	□	■	■	□	■
ExpansionHunter		■			■				■	■		■
GATK <sup>b</sup>		■		□		□			■			
LUMPY					□	■				□		□
CNVnator					□				■			■
RUFUS		■								■		
CNVkit								■				■
BreakDancer					□	■						
Illumina DRAGEN depth-based CNV caller	■											
SvABA: SV/indel Analysis by Assembly				■								
CoNIFER <sup>c</sup>							■					
ERDS: estimation by reads depth w/ SNVs									■			
BreakSeq2					□							
DELLY2					□							
<b>Jointly call and/or genotype SVs</b>												
smoove							■			■		■
SVTyper					□					□		□
<b>Annotate SVs</b>												
AnnotSV		■					■	■		■	■	■
gnomAD-SV					■							■
duphold										□		□
<b>Run or combine output from other tools</b>												
XHMM		■		■		■						
SURVIVOR					□							■
Parliament2					■							

■ Tool called directly. □ Tool called indirectly (e.g., by a wrapper).

Each SV calling tool identifies subsets of SVs by type or other factors, and so in practice, the output of multiple methods must typically be combined and considered together. Wrapper tools that automatically call and combine results from multiple other SV detection methods improve the efficiency of this process. Duke/Columbia, NIH, Stanford, and Vanderbilt only use SV calling tools in specific cases or contexts rather than as part of their regular pipelines. Tool citations are listed in Extended Data Table 1.

CNV copy-number variant, SNV single-nucleotide variant.

<sup>a</sup>Manta is used by BaylorSeq to generate putative SV calls, which are then shared with the clinical sites.

<sup>b</sup>The two functions from GATK used are GermlineCNVCaller and DepthOfCoverage (DoC); the latter is used to detect exonic deletions or duplications.

<sup>c</sup>In contrast to other tools, CoNIFER runs on exome sequencing (ES) data rather than genome sequencing (GS) data.

automated annotation and filtering of these variants is the next step of the variant prioritization process (Fig. 1i, Extended Data Table 2). Annotations fall into four categories: (1) known disease associations, (2) prevalence across healthy human populations, (3) predicted pathogenicity and functional effect, and (4) inheritance. Many scores exist across the first three categories;<sup>25</sup> in the following sections we explore those that are used in practice for rare disease diagnosis.

#### Known disease-associated genes

Many specific genetic variants have previously been determined to cause human disease, and it is useful to first look for the presence of these variants in a patient's sequencing data. Databases compiling disease-causing variants, the genes they impact, and their phenotypic associations are used by ten clinical sites (Table 3). Genetic testing laboratories, including BaylorSeq,

use these in addition to internal databases containing similar information. Disease-relevant variants are listed on clinical reports and are considered during the initial pass of each UDN case at all clinical sites.

#### Variant segregation in healthy human populations

Several positions within the human genome naturally vary across healthy individuals, and "common" variants at these positions are unlikely to cause the conditions under investigation by the UDN. Though rare combinations of otherwise common variants may lead to disease,<sup>26</sup> clinical sites do not currently consider all common variant combinations. Instead, variants observed more than 1 in 100 times across healthy populations (i.e., minor allele frequency [MAF] > 0.01) are typically excluded during the first pass of the data. The exact MAF threshold used depends on the suspected mode of inheritance. Lower MAF thresholds are used

**Table 2.** Quality control (QC) checks of variants for rare disease diagnosis.

	BCFtools <sup>a</sup>	vcf.lobio	GATK DoC	IGV	Peddy	PCA	LAMP-LD	plink	SNP chips	ATAV	denovo-db <sup>b</sup>	novoCaller <sup>b</sup>	Salvage Pathway <sup>b</sup>
<b>Sequencing quality is acceptable</b>													
Density of variants is uniform genome-wide and for each chromosome	•••••	•											
Transition/transversion ratio is ~2 for ES or ~3 for GS <sup>c</sup>	•••••	•											
Homozygous/heterozygous site ratio in ~1.5	•••••												
Most variants are SNVs but all variant types are present	•••••	•											
Distributions of variant read depth, variant quality and genotype quality have no unexpected patterns or outliers	•••••	•		•									
Variants are present in aligned reads										⑪			
<b>Sequenced samples match expectation</b>													
Sex matches expectation				•	•				•	•			
Ancestry (relative to 1000 Genomes samples) matches expectation					•	•	•	•					
Relatedness between individuals matches expectation <sup>d</sup>					••					•			
<b>Low rate of Mendelian violations</b>													
Mosaicism (i.e., 2+ nucleotides not in a 50-50 ratio at multiple sites) is not present unless expected <sup>e</sup>									•	•			
Regions that match between related individuals occur in contiguous stretches								•					
Sites that are homozygous in child are at least heterozygous in both parents						⑪		•		•			
Low <i>de novo</i> variant count (<10 for ES and <75 for GS)								•			•	•	•
Number of clinical sites using tool for QC step: • one •• two ••• three •••• four ••••• five ⑪ eleven													

QC checks of variant data fall into three main categories, listed in bold above. Although some tools can be used for many of these steps, we illustrate here which QC steps they are actually used for in practice. Note the clarifications for some of the QC tools and steps listed in footnotes a–e. Tool citations are listed in Extended Data Table 1.

ES exome sequencing, GS genome sequencing, SNV single-nucleotide variant.

<sup>a</sup>BCFtools refers to the Wellcome Trust Sanger Institute's suite of tools: BCFtools, VCFtools, SAMtools, and HTSLib.

<sup>b</sup>These tools either call *de novo* variants from sequencing reads to reduce false positive calls or provide *de novo* frequencies where a high frequency indicates a likely false positive.

<sup>c</sup>The expected transition (Ts) to transversion (Tv) ratios assume variants are called with respect to the human reference sequence; if variants are called with respect to computed ancestral alleles, the expected Ts/Tv ratio for ES should be ~1.

<sup>d</sup>Expected relatedness between family members is estimated using a "kinship coefficient"; unexpectedly low kinship implies a family member is not as related as was originally assumed, unexpectedly high kinship suggests consanguinity, and maximal kinship implies an accidental sample duplication.

<sup>e</sup>Mosaicism—where an individual contains a mix of genetically distinct cells—may be relevant for disease rather than only indicative of sequencing errors.

for suspected dominant conditions because the variants causing the extremely rare phenotypes of UDN patients are assumed to be naturally selected against and thus equally rare in the general population and entirely absent in control population databases. Higher MAF thresholds are used for suspected recessive conditions because heterozygous individuals would not be expected to manifest severe disease features.

All UDN sites use data from the Broad Institute's Genome Aggregation Database (gnomAD) to compute MAFs, and seven sites also compute MAFs from smaller or population-specific data sets on a case-by-case basis (Table 3). Two sites eliminate variants that are homozygous in three or more healthy individuals in these

data sets. At the NIH site, rather than thresholding on MAFs computed directly from variant proportions in gnomAD, 95% Wilson confidence score intervals computed from these proportions are used to retain rare variants occurring in low coverage regions. Finally, five sites flag variants that are present in data sets internal to their institutions, because variants present in asymptomatic or differently symptomatic individuals are unlikely to be disease-relevant.

Eight sites consult SV databases to check the existence and/or MAF of detected SVs (Table 3, Extended Data Table 1). Multiple databases are checked in practice because the SV detection tools used across databases differ, so the absence or rarity of an SV in

**Table 3.** Human genetic variation data sets and derived tools.

	BaylorSeq	BCM	Duke/Columbia	Harvard	Miami	NIH	PacificNW	Stanford	UCLA	Utah	Vanderbilt	WUSTL
<b>Known disease gene databases</b>												
ClinVar	●	●	●	●	●	●	●	●	●	●	●	●
OMIM	●	●	●	●	●	●	●	●	●	●	●	●
HGMD: Human Gene Mutation Database	●	●	●	●	●	●	●	●	●	●	●	●
dbSNP	●	●	●				●			●		
CGD: Clinical Genomic Database								●		●	●	
Orphanet								●		●		
<b>Healthy human population single-nucleotide variant (SNV)/indel databases</b>												
gnomAD: Genome Aggregation Database	●	●	●	●	●	●	●	●	●	●	●	●
ExAC: Exome Aggregation Consortium	●	●	●	●	●	●	●	●	●	●	○	●
1000 Genomes Project	●	●	●	●	●	●	●	●	●	●	●	●
Institution—internal controls <sup>a</sup>		●	●	●		●	●	●	●	●	●	
EVS: Exome Variant Server	●		●	●	●		●	●	●	●	●	
TOPMed: Trans-Omics for Precision Medicine			○				●	●	●	○		○
UK10K							●	●	●			○
Greater Middle East (GME) Variome Project			○									
xKJPN: 1000+ Japanese			○									
GenomeAsia 100 K Project			○									
Iranome			○									
<b>Human structural variant (SV) databases</b>												
gnomAD-SV: Genome Aggregation Database SVs	●	●	●	○	●	●	●	●	●	●	●	●
DGV: Database of Genomic Variants	●	●	●	○		●	●	●	●	●	●	●
dbVar: Database of Genomic Structural Variation	●	●	●	○			●	●	●	●	●	●
ClinGen: Clinical Genome Resource	●	●	●	○			●	●	●	●	●	●
DECIPHER	●	●	●	○			●	●	●	●	●	●
Institution—internal controls <sup>a</sup>				○					●	●	●	●
<b>Within-human selective constraint scores</b>												
pLI: probability of loss-of-function (LoF) intolerance	●	●	●	●	●	●	●	○	●	●	●	●
Missense (constraint) Z score	●	●	●	○	●	●	●	●	●	●	●	●
pREC: probability of homozygote LoF intolerance (sub)RVIS: Residual Variation Intolerance Score	●	●	●	○	●	●	●	●	●	●	●	●
L-o/e-UF: LoF observed/expected upper-bound fraction	●	●	●	●	●	●	●	●	●	●	●	●
CCR: constrained coding regions				●			●	●	●	●	●	●

**Table 3** continued

	BaylorSeq	BCM	Duke/Columbia	Harvard	Miami	NIH	PacificNW	Stanford	UCLA	Utah	Vanderbilt	WUSTL
LIMBR: Localized Intolerance Model w/ Bayesian Regression		●										
MTR: missense tolerance ratio		●										
s_het: selective effect of heterozygous LoF				●								
M-o/e-UF: missense observed/expected upper-bound fraction						●						
LoFtool											●	
● Tool used by default. ○ Tool used in specific cases or contexts only. <sup>b</sup>												

Knowledge of variation within human populations with and without disease can be effectively used to assess the likelihood of a variant to cause the genetic condition under investigation. Tool and data set citations are listed in Extended Data Table 1.

<sup>a</sup>Human sequence variation data sets that are internal to particular institutions and used by clinical sites surveyed here include variants present in patients from Baylor College of Medicine (BCM), the Institute for Genomic Medicine (Duke/Columbia), Brigham Genomic Medicine (Harvard), the NIH Undiagnosed Diseases Program (NIH), Centers for Mendelian Genomics (PacificNW), University of California–Los Angeles (UCLA), the Centre d'Etude du Polymorphisme Humain (Utah), and BioVu (Vanderbilt), and a curated set of copy-number variants (CNVs) detected via genome sequencing (GS) and confirmed via chromosomal microarray analysis (Washington University School of Medicine [WUSTL]).

<sup>b</sup>The contexts in which specific human population variant data sets are used include historical reasons (ExAC), when a variant's gnomAD-derived MAF is 0 or close to 0 (TOPMed), when patients' inferred ancestry is non-European (TOPMed), Middle Eastern (GME), Japanese (xKJPN), Asian (GenomeAsia), and/or Iranian (Iranome), and when a predicted structural variant impacts a clinically relevant gene (gnomAD-SV, DGV, ClinGen, DECIPHER).

one database may reflect a particular SV detection approach rather than true population rarity.

Simple genetic variation observed across healthy humans tends to be sparsely distributed with varying degrees of impact. These features can be used to capture how regions of the human genome may be intolerant of loss-of-function (LoF) variants, such as frameshift or protein-truncating variants. Nine surveyed sites incorporate selective constraint scores derived from and released with gnomAD data in their diagnostic pipelines, with the probability of heterozygous LoF intolerance scores and missense constraint Z scores used most commonly (Table 3).

#### Predicted pathogenicity and functional effect of variants

Various tools predict the pathogenicity of uncovered variants.<sup>25</sup> Values derived from cross-species comparative genomics contribute heavily to pathogenicity predictors, as positions that are conserved across species tend to be functionally critical. However, since most candidate coding variants are evolutionarily well-conserved, only five sites directly consider conservation in their diagnostic pipelines (Table 4, Extended Data Table 1).

The most commonly used pathogenicity predictors for rare disease diagnosis—used by eight clinical sites each—are Combined Annotation Dependent Depletion (CADD) and Rare Exome Variant Ensemble Learner (REVEL), each of which consider multiple variant annotations and where scores >25 and >0.3 respectively indicate likely pathogenic variants. Nearly all predicted pathogenicity scores used, with the exception of ReMM, indicate disease relevance primarily for coding variants.<sup>27</sup>

Indeed, predicting and experimentally validating the pathogenic impact of noncoding variants is notoriously difficult. All 12 sites use tools to predict how noncoding variants alter expected gene expression and splicing. Few sites use the same subset of tools for this task, though SpliceAI is the most commonly used tool overall (Table 4).

#### Mode of inheritance

After variants have been quality checked, MAF filtered, and annotated, Mendelian mode of inheritance is evaluated next by the clinical sites. Some sites simultaneously consider the functional impact of variants, where, for instance, intergenic or perceived synonymous variants are excluded.<sup>3</sup> Despite the ubiquity of this step, each site uses different tools for computing inheritance patterns.

For a dominantly inherited genetic condition to manifest, only one defective copy of the relevant gene is required, whereas recessive disease manifestation requires two defective gene copies. GS of unrelated or distantly related affected individuals is desired in suspected dominant cases to find rare, shared variants.

In sporadic cases—caused by a single *de novo* dominant or two recessive variants—GS of at least the affected individual and both unaffected parents is desired. Selecting heterozygous variants in the affected individual that are absent in both unaffected parents or homozygous variants in the affected individual that are absent in at least one parent via straightforward segregation analysis results in a majority of spurious *de novo* calls. These false positive calls stem from inadequate sequence coverage or alignment in parents from whom variants were in fact inherited and/or inaccurate modeling of underlying variant frequencies. Four sites regularly use specialized *de novo* calling tools or databases to offset these issues (Table 2). Fixing *de novo* calling errors requires analysis of sequencing reads, which many genetic testing centers do not readily provide.

Occasionally in sporadic and/or recessive cases, the same disease-causing variant is inherited from both heterozygous parents and can be easily detected as a homozygous variant. Genomic regions containing only homozygous variants in an affected individual with

**Table 4.** Tools for assigning the pathogenic likelihood or functional impact of variants.

	BaylorSeq	BCM	Duke/Columbia	Harvard	Miami	NIH	PacificNW	Stanford	UCLA	Utah	Vanderbilt	WUSTL
<b>Cross-species conservation scores</b>												
GERP++: Genomic Evolutionary Rate Profiling	●	●	●			●	●				●	
PhastCons				●	●	●	●					
<b>Predicted functionality or pathogenicity</b>												
PolyPhen-2	●	●	●	●	●	●	●	●	●	●		
SIFT	●	●		●	●	●	●	●	●	●		
MutationTaster			●	●	●	●	●	●				
MVP: missense variant pathogenicity								●				
ReMM: regulatory Mendelian mutation								●				
<b>Ensemble pathogenicity predictors</b>												
CADD: Combined Annotation Dependent Depletion	●	●			●	●	●	●	●		●	●
REVEL: Rare Exome Variant Ensemble Learner		●		●		●	●	●		●		●
DANN: Deep Neural Net version of CADD	●		●		●							●
M-CAP: Mendelian Clinically Applicable Pathogenicity								●				
DOMINO: Dominant Disorder Associated Genes <sup>a</sup>		●										
Eigen											●	
<b>Predicted splice- or expression-altering effect</b>												
SpliceAI	●	●		●	●		●	●	●	●	●	●
GTEX: Genotype-Tissue Expression		●			●			●		●		
SpliceRegion annotations from VEP							●	●		●		●
dbSCSNV (splicing consensus SNVs)								●	●			●
Human Splicing Factor					●						●	
MMSplice: Modular modeling of splicing		●			●							●
MaxEntScan											●	
Trap: Transcript-inferred Pathogenicity			●									

Variants of uncertain significance (i.e., that are not already known to be associated with disease) can be evaluated for functional or pathogenic impact using predictive models. Tool citations are listed in Extended Data Table 1.

<sup>a</sup>Unlike other tools, DOMINO provides scores per gene rather than per variant.



nonconsanguineous parents can also indicate an inherited deletion from one parent or uniparental isodisomy. These latter phenomena, revealed as Mendelian violations during the QC process (Table 2), can manifest in a recessive disease despite only one parent being heterozygous for the disease-causing variant. Often in undiagnosed recessive cases, two or more different heterozygous variants, each either inherited or occurring *de novo*, can give rise to the disease phenotype; these variants are referred to as compound heterozygous pairs. The complete set of compound heterozygous variant pairs in any given case is very large, and so filters—such as restricting to rare, LoF, likely pathogenic variants—are applied beforehand. If too few candidate explanatory variants pass these filters, the NIH, WUSTL and Miami sites use internal “second tier” schemes, such as increasing the allowable MAF threshold, to rescue additional compound heterozygous pairs.<sup>28</sup>

#### Integration of nonsequencing data

Cases with nondiagnostic genetic testing have eventually been solved by reanalysis approaches that leverage additional data, such as transcriptome sequencing<sup>29,30</sup> (RNA-seq) or “deep phenotyping,”<sup>31,32</sup> to complement ES and GS.

#### Transcriptome sequencing

RNA-seq is increasingly utilized to (1) confirm suspected expression- or splice-altering variants initially prioritized through genomic sequencing, and/or (2) highlight genes that are aberrantly expressed relative to healthy, tissue-matched samples from databases such as GTEx (<https://gtexportal.org/>).<sup>29,30</sup> BCM, Stanford, and UCLA regularly use RNA-seq data for variant prioritization, and two other sites are actively working to incorporate RNA-seq data into their workflows as well (Extended Data Table 3). Vanderbilt uses PrediXcan to correlate observed phenotypes with imputed, rather than directly measured, gene expression.<sup>33</sup>

#### Structured phenotyping

Deep phenotyping of patients is critical to the overall UDN process (Fig. 1a) and enables clinicians to focus on genes associated with a patient’s symptoms or suspected disease. Symptom terms are standardized via the Human Phenotype Ontology (HPO) and explicitly annotated for each UDN case during the in-person evaluation.<sup>34</sup> Computational tools can reason over these terms to generate gene panels that complement manual efforts.<sup>35</sup> All clinical sites have access to genes ranked by PhenoTips, a program embedded into the UDN data server. Eight clinical sites and BaylorSeq use additional tools to prioritize genes from patients’ phenotypes (Fig. 1j, Extended Data Table 4).<sup>36</sup> Amelie is used by five sites to scour the literature for examples of genes causing patients’ observed phenotypes, a process typically performed manually using the Monarch Initiative’s gene–phenotype browser. Exomiser is used by three sites to integrate genotype–phenotype data and runs in parallel to existing pipelines. Finally, pairwise associations between genes and HPO terms are downloadable from the HPO website; the union of genes associated with all annotated HPO terms per patient can be used directly or intersected with sets of disease-relevant genes from OMIM and HGMD. This approach is used by three sites regularly but has been implemented for various projects at all clinical sites.

#### Workflow management and wrapper tools

The complex workflows described here must be well-documented, customizable per case, and provide results in a timely manner and intuitive format. Case materials should be accessible by collaborative teams of clinicians, bioinformaticians, and genetic counselors. In practice, all sites use automated platforms to call, annotate, and prioritize candidate diagnostic variants (Extended Data Table 5, Extended Data Table 6). Spreadsheets are the most common tool

used by all sites for storing, sharing, and commenting on variant-level data. Many sites also use commercial solutions for case management, which has enabled secure transition of certain workflow components to the cloud.

## DISCUSSION

Pinpointing the genetic variants giving rise to ultrarare, undiagnosed diseases is a challenging and pressing problem being tackled on a case-by-case basis by clinical researchers worldwide. The computational tools utilized during these investigative efforts reflect relevant community standards but can also diverge across institutions and even across cases handled by the same clinical team.

The diverse, exploratory techniques employed by UDN clinical sites can overcome inherent limitations of clinical case review and standard sequencing interpretation provided by genetic testing laboratories—both of which rely on existing disease gene knowledge—by uncovering novel disease loci. For instance, when no compelling variants were found in phenotypically prioritized genes in two patients presenting with muscular and white matter abnormalities, a genetics-driven UDN pipeline uncovered diagnostic *de novo* missense variants in both individuals in *TOMM70*, a gene previously unassociated with disease.<sup>37</sup> Similarly, sequencing analyses were able to uncover *de novo*, heterozygous variants in nine individuals with neurodevelopmental delay and other multisystem anomalies in *CDH2*, a gene previously unassociated with a Mendelian neurodevelopmental condition.<sup>38</sup>

Indeed, divergent aspects of UDN pipelines reflect promising avenues for case reanalysis and reveal areas where technical developments would be most impactful. Improving SV detection specificity would aid in cases with nondiagnostic microarrays, gene panels, and GS. Experimentally verifiable pathogenicity predictions for noncoding variants may solve cases with nondiagnostic ES. Finally, automated integration of additional data, such as RNA-seq,<sup>29,30</sup> long-read sequencing,<sup>39</sup> and epigenetic modifications,<sup>40</sup> may also increase the diagnostic rate for cases with inconclusive GS.

Consensus tools used across sites by multiple clinical research teams have been convincingly evaluated and are easily incorporated into existing workflows external to their original development environment. Clinical sites strive to incorporate better tools—including those developed in-house—as they emerge over time. Flexible, open-source implementations ease this process and can ultimately shorten the time to and improve the rate of diagnosis. Initiatives like the UDN provide an excellent opportunity to assess and share tools and ideas and jointly develop methods inspired by the most challenging undiagnosed cases.

## DATA AVAILABILITY

All data used in this analysis are available in the Main and Extended Data Tables.

Received: 11 August 2020; Revised: 14 December 2020; Accepted: 17 December 2020;

Published online: 12 February 2021

## REFERENCES

1. Boycott, K. M., Vanstone, M. R., Bulman, D. E. & MacKenzie, A. E. Rare-disease genetics in the era of next-generation sequencing: discovery to translation. *Nat. Rev. Genet.* **14**, 681–691 (2013).
2. Online Mendelian Inheritance in Man, OMIM. (McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University, Baltimore, MD). <https://omim.org>.
3. Robinson, P. N. et al. Improved exome prioritization of disease genes through cross-species phenotype comparison. *Genome Res.* **24**, 340–348 (2014).

4. Karczewski, K. J. et al. The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **581**, 434–443 (2020).
5. Adzhubei, I. A. et al. A method and server for predicting damaging missense mutations. *Nat. Methods* **7**, 248–249 (2010).
6. Posey, J. E. et al. Insights into genetics, human biology and disease gleaned from family based genomic studies. *Genet. Med.* **21**, 798–812 (2019).
7. Splinter, K. et al. Effect of genetic diagnosis on patients with previously undiagnosed disease. *N. Engl. J. Med.* **379**, 2131–2139 (2018).
8. Macnamara, E. F. et al. Cases from the Undiagnosed Diseases Network: The continued value of counseling skills in a new genomic era. *J. Genet. Couns.* **28**, 194–201 (2019).
9. Macnamara, E. F. & D'Souza, P, Undiagnosed Diseases Network & Tiffit, C. J. The undiagnosed diseases program: approach to diagnosis. *Transl. Sci. Rare Dis.* **4**, 179–188 (2020).
10. Wambach, J. A. et al. Functional characterization of biallelic RTTN variants identified in an infant with microcephaly, simplified gyral pattern, pontocerebellar hypoplasia, and seizures. *Pediatr. Res.* **84**, 435–441 (2018).
11. Lee, H. et al. Clinical exome sequencing for genetic identification of rare Mendelian disorders. *JAMA* **312**, 1880–1887 (2014).
12. Haghghi, A. et al. An integrated clinical program and crowdsourcing strategy for genomic sequencing and Mendelian disease gene discovery. *NPJ Genom. Med.* **3**, 21 (2018).
13. Richards, S. et al. Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genet. Med.* **17**, 405–424 (2015).
14. Philippakis, A. A. et al. The Matchmaker Exchange: a platform for rare disease gene discovery. *Hum. Mutat.* **36**, 915–921 (2015).
15. Frost, J. H. & Massagli, M. P. Social uses of personal health information within PatientsLikeMe, an online patient community: what can happen when patients have access to one another's data. *J. Med. Internet Res.* **10**, e15 (2008).
16. Wang, J. et al. MARRVEL: integration of human and model organism genetic resources to facilitate functional annotation of the human genome. *Am. J. Hum. Genet.* **100**, 843–853 (2017).
17. Bimber, B. N., Yan, M. Y., Peterson, S. M. & Ferguson, B. mGAP: the macaque genotype and phenotype resource, a framework for accessing and interpreting macaque variant data, and identifying new models of human disease. *BMC Genomics* **20**, 176 (2019).
18. Meyer, E. et al. Mutations in the histone methyltransferase gene KMT2B cause complex early-onset dystonia. *Nat. Genet.* **49**, 223–237 (2017).
19. Regier, A. A. et al. Functional equivalence of genome sequencing analysis pipelines enables harmonized variant calling across human genetics projects. *Nat. Commun.* **9**, 4038 (2018).
20. Van der Auwera, G. A. et al. From FastQ data to high confidence variant calls: the Genome Analysis Toolkit best practices pipeline. *Curr. Protoc. Bioinformatics* **43**, 11.10.1–11.10.33 (2013).
21. Poplin, R. et al. A universal SNP and small-indel variant caller using deep neural networks. *Nat. Biotechnol.* **36**, 983–987 (2018).
22. Collins, R. L. et al. A structural variation reference for medical and population genetics. *Nature* **581**, 444–451 (2020).
23. Mahmoud, M. et al. Structural variant calling: the long and the short of it. *Genome Biol.* **20**, 246 (2019).
24. Kosugi, S. et al. Comprehensive evaluation of structural variation detection algorithms for whole genome sequencing. *Genome Biol.* **20**, 117 (2019).
25. Liu, X., Wu, C., Li, C. & Boerwinkle, E. dbNSFP v3.0: a one-stop database of functional predictions and annotations for human nonsynonymous and splice-site SNVs. *Hum. Mutat.* **37**, 235–241 (2016).
26. Posey, J. E. Genome sequencing and implications for rare disorders. *Orphanet J. Rare Dis.* **14**, 153 (2019).
27. Mather, C. A. et al. CADD score has limited clinical validity for the identification of pathogenic variants in noncoding regions in a hereditary cancer panel. *Genet. Med.* **18**, 1269–1275 (2016).
28. Gu, F. et al. A suite of automated sequence analyses reduces the number of candidate deleterious variants and reveals a difference between probands and unaffected siblings. *Genet. Med.* **21**, 1772–1780 (2019).
29. Lee, H. et al. Diagnostic utility of transcriptome sequencing for rare Mendelian diseases. *Genet. Med.* **22**, 490–499 (2020).
30. Frésard, L. et al. Identification of rare-disease genes using blood transcriptome sequencing and large control cohorts. *Nat. Med.* **25**, 911–919 (2019).
31. Shashi, V. et al. A comprehensive iterative approach is highly effective in diagnosing individuals who are exome negative. *Genet. Med.* **21**, 161–172 (2019).
32. Pena, L. D. M. et al. Looking beyond the exome: a phenotype-first approach to molecular diagnostic resolution in rare and undiagnosed diseases. *Genet. Med.* **20**, 464–469 (2018).
33. Gamazon, E. R. et al. A gene-based association method for mapping traits using reference transcriptome data. *Nat. Genet.* **47**, 1091–1098 (2015).
34. Köhler, S. et al. Expansion of the Human Phenotype Ontology (HPO) knowledge base and resources. *Nucleic Acids Res.* **47**, D1018–D1027 (2019).
35. Smedley, D. & Robinson, P. N. Phenotype-driven strategies for exome prioritization of human Mendelian disease genes. *Genome Med.* **7**, 81 (2015).
36. Gonzalez, M. et al. Innovative genomic collaboration using the GENESIS (GEM. app) platform. *Hum. Mutat.* **36**, 950–956 (2015).
37. Dutta, D. et al. De novo mutations in TOMM70, a receptor of the mitochondrial import translocase, cause neurological impairment. *Hum. Mol. Genet.* **29**, 1568–1579 (2020).
38. Accogli, A. et al. De novo pathogenic variants in N-cadherin cause a syndromic neurodevelopmental disorder with corpus callosum, axon, cardiac, ocular, and genital defects. *Am. J. Hum. Genet.* **105**, 854–868 (2019).
39. Merker, J. D. et al. Long-read genome sequencing identifies causal structural variation in a Mendelian disease. *Genet. Med.* **20**, 159–163 (2017).
40. Turro, E. et al. Whole-genome sequencing of patients with rare diseases in a national health system. *Nature* **583**, 96–102 (2020).

## ACKNOWLEDGEMENTS

Thank you to the UDN Tool Building Coalition for discussions about tools in use or under development, to Daniel Traviglia for clarifications on UDN data availability, and to Rebecca Reimers for writing feedback. Research reported here was supported by the NIH Common Fund, through the Office of Strategic Coordination/Office of the NIH Director under award numbers U01HG007530, U01HG007942, U01HG007672, U01HG007690, U01HG010218, U01HG007703, U01HG010230, U01HG010217, U01HG010233, U01HG007674, and U01HG010215, and by the Intramural Research Program of the National Human Genome Research Institute. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

## AUTHOR CONTRIBUTIONS

Conceptualization: S.N.K., S.R.S., I.S.K. Data curation: S.N.K., D.B., M.V., J.B.K., B.N.P., S.Z., E.B., H.L., A.H., L.B., A.B., J.C., S.M., A.A., D.R.M., P.L., D.J.W., A.J.P. Formal analysis: S.N.K. Funding acquisition: I.S.K. Investigation: S.N.K., S.R.S., I.S.K. Methodology: S.N.K. Visualization: S.N.K.; Writing—original draft: S.N.K. Writing—review & editing: S.N.K., D.B., M.V., K.L., C.E., S.R.S.

## COMPETING INTERESTS

P.L. is an employee of Baylor College of Medicine and derives support through a professional services agreement with Baylor Genetics, which performs clinical genetic testing services. The other authors declare no competing interests.

## ADDITIONAL INFORMATION

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41436-020-01084-8>.

**Correspondence** and requests for materials should be addressed to I.S.K.

**Reprints and permission information** is available at <http://www.nature.com/reprints>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

**UNDIAGNOSED DISEASES NETWORK**

Maria T. Acosta<sup>16</sup>, Margaret Adam<sup>17</sup>, David R. Adams<sup>16</sup>, Pankaj B. Agrawal<sup>18</sup>, Mercedes E. Alejandro<sup>19</sup>, Justin Alvey<sup>20</sup>, Laura Amendola<sup>17</sup>, Ashley Andrews<sup>20</sup>, Euan A. Ashley<sup>21</sup>, Mahshid S. Azamian<sup>19</sup>, Carlos A. Bacino<sup>19</sup>, Guney Bademci<sup>22</sup>, Eva Baker<sup>16</sup>, Ashok Balasubramanyam<sup>19</sup>, Dustin Baldrige<sup>23,24</sup>, Jim Bale<sup>20</sup>, Michael Bamshad<sup>17</sup>, Deborah Barbooth<sup>22</sup>, Pinar Bayrak-Toydemir<sup>20</sup>, Anita Beck<sup>17</sup>, Alan H. Beggs<sup>18</sup>, Edward Behrens<sup>25</sup>, Gill Bejerano<sup>21</sup>, Jimmy Bennett<sup>17</sup>, Beverly Berg-Rood<sup>17</sup>, Jonathan A. Bernstein<sup>21</sup>, Gerard T. Berry<sup>18</sup>, Anna Bican<sup>26</sup>, Stephanie Bivona<sup>22</sup>, Elizabeth Blue<sup>17</sup>, John Bohnsack<sup>20</sup>, Carsten Bonnenmann<sup>16</sup>, Devon Bonner<sup>21</sup>, Lorenzo Botto<sup>20</sup>, Brenna Boyd<sup>17</sup>, Lauren C. Briere<sup>18</sup>, Elly Brokamp<sup>26</sup>, Gabrielle Brown<sup>27</sup>, Elizabeth A. Burke<sup>16</sup>, Lindsay C. Burrage<sup>19</sup>, Manish J. Butte<sup>27</sup>, Peter Byers<sup>17</sup>, William E. Byrd<sup>28</sup>, John Carey<sup>20</sup>, Olveen Carrasquillo<sup>22</sup>, Ta Chen Peter Chang<sup>22</sup>, Sirisak Chanprasert<sup>17</sup>, Hsiao-Tuan Chao<sup>19</sup>, Gary D. Clark<sup>19</sup>, Terra R. Coakley<sup>21</sup>, Laurel A. Cobban<sup>18</sup>, Joy D. Cogan<sup>26</sup>, Matthew Coggins<sup>18</sup>, F. Sessions Cole<sup>23</sup>, Heather A. Colley<sup>16</sup>, Cynthia M. Cooper<sup>18</sup>, Heidi Cope<sup>29</sup>, William J. Craigen<sup>19</sup>, Andrew B. Crouse<sup>28</sup>, Michael Cunningham<sup>17</sup>, Precilla D'Souza<sup>16</sup>, Hongzheng Dai<sup>19</sup>, Surendra Dasari<sup>30</sup>, Joie Davis<sup>16</sup>, Jyoti G. Daya<sup>1</sup>, Matthew Deardorff<sup>25</sup>, Esteban C. Dell'Angelica<sup>27</sup>, Shweta U. Dhar<sup>19</sup>, Katrina Dipple<sup>17</sup>, Daniel Doherty<sup>17</sup>, Naghmeh Dorrani<sup>27</sup>, Argenia L. Doss<sup>16</sup>, Emilie D. Douine<sup>27</sup>, David D. Draper<sup>16</sup>, Laura Duncan<sup>26</sup>, Dawn Earl<sup>17</sup>, David J. Eckstein<sup>16</sup>, Lisa T. Emrick<sup>19</sup>, Christine M. Eng<sup>31</sup>, Cecilia Esteves<sup>32</sup>, Marni Falk<sup>25</sup>, Liliana Fernandez<sup>21</sup>, Carlos Ferreira<sup>16</sup>, Elizabeth L. Fieg<sup>18</sup>, Laurie C. Findley<sup>16</sup>, Paul G. Fisher<sup>21</sup>, Brent L. Fogel<sup>27</sup>, Irman Forghani<sup>22</sup>, Laure Fresard<sup>21</sup>, William A. Gahl<sup>16</sup>, Ian Glass<sup>17</sup>, Bernadette Gochuico<sup>16</sup>, Rena A. Godfrey<sup>16</sup>, Katie Golden-Grant<sup>17</sup>, Alica M. Goldman<sup>19</sup>, Madison P. Goldrich<sup>16</sup>, David B. Goldstein<sup>33</sup>, Alana Grajewski<sup>22</sup>, Catherine A. Groden<sup>16</sup>, Irma Gutierrez<sup>27</sup>, Sihoun Hahn<sup>17</sup>, Rizwan Hamid<sup>26</sup>, Neil A. Hanchard<sup>19</sup>, Kelly Hassey<sup>25</sup>, Nichole Hayes<sup>23</sup>, Frances High<sup>18</sup>, Anne Hing<sup>17</sup>, Fuki M. Hisama<sup>17</sup>, Ingrid A. Holm<sup>18</sup>, Jason Hom<sup>21</sup>, Martha Horike-Pyne<sup>17</sup>, Alden Huang<sup>27</sup>, Yong Huang<sup>21</sup>, Laryssa Huryn<sup>16</sup>, Rosario Isasi<sup>22</sup>, Fariha Jamal<sup>19</sup>, Gail P. Jarvik<sup>17</sup>, Jeffrey Jarvik<sup>17</sup>, Suman Jayadev<sup>17</sup>, Lefkothea Karaviti<sup>19</sup>, Jennifer Kennedy<sup>26</sup>, Dana Kiley<sup>23</sup>, Isaac S. Kohane<sup>32</sup>, Jennefer N. Kohler<sup>21</sup>, Susan Korrick<sup>18</sup>, Mary Kozuira<sup>26</sup>, Deborah Krakow<sup>27</sup>, Donna M. Krasnewich<sup>16</sup>, Elijah Kravets<sup>21</sup>, Joel B. Krier<sup>18</sup>, Grace L. LaMoure<sup>16</sup>, Seema R. Lalani<sup>19</sup>, Byron Lam<sup>22</sup>, Christina Lam<sup>17</sup>, Brendan C. Lanpher<sup>30</sup>, Ian R. Lanza<sup>30</sup>, Lea Latham<sup>16</sup>, Kimberly LeBlanc<sup>32</sup>, Brendan H. Lee<sup>19</sup>, Hane Lee<sup>27</sup>, Roy Levitt<sup>22</sup>, Richard A. Lewis<sup>19</sup>, Sharyn A. Lincoln<sup>18</sup>, Pengfei Liu<sup>31</sup>, Xue Zhong Liu<sup>22</sup>, Nicola Longo<sup>20</sup>, Sandra K. Loo<sup>27</sup>, Joseph Loscalzo<sup>18</sup>, Richard L. Maas<sup>18</sup>, John MacDowall<sup>16</sup>, Calum A. MacRae<sup>18</sup>, Ellen F. Macnamara<sup>16</sup>, Valerie V. Maduro<sup>16</sup>, Marta M. Majcherska<sup>21</sup>, Bryan C. Mak<sup>27</sup>, May Christine V. Malicdan<sup>16</sup>, Laura A. Mamounas<sup>16</sup>, Teri A. Manolio<sup>16</sup>, Rong Mao<sup>20</sup>, Kenneth Maravilla<sup>17</sup>, Thomas C. Markello<sup>16</sup>, Ronit Marom<sup>19</sup>, Gabor Marth<sup>20</sup>, Beth A. Martin<sup>21</sup>, Martin G. Martin<sup>27</sup>, Julian A. Martinez-Agosto<sup>27</sup>, Shruti Marwaha<sup>21</sup>, Jacob McCauley<sup>22</sup>, Allyn McConkie-Rosell<sup>29</sup>, Colleen E. McCormack<sup>21</sup>, Alexa T. McCray<sup>32</sup>, Elisabeth McGee<sup>27</sup>, Heather Mefford<sup>17</sup>, J. Lawrence Merritt<sup>17</sup>, Matthew Might<sup>28</sup>, Ghayda Mirzaa<sup>17</sup>, Eva Morava<sup>30</sup>, Paolo M. Moretti<sup>19</sup>, Paolo Moretti<sup>20</sup>, Deborah Mosbrook-Davis<sup>16</sup>, John J. Mulvihill<sup>16</sup>, David R. Murdock<sup>19</sup>, Anna Nagy<sup>32</sup>, Mariko Nakano-Okuno<sup>28</sup>, Avi Nath<sup>16</sup>, Stanley F. Nelson<sup>27</sup>, John H. Newman<sup>26</sup>, Sarah K. Nicholas<sup>19</sup>, Deborah Nickerson<sup>17</sup>, Shirley Nieves-Rodriguez<sup>27</sup>, Donna Novacic<sup>16</sup>, Devin Oglesbee<sup>30</sup>, James P. Orenge<sup>19</sup>, Laura Pace<sup>20</sup>, Stephen Pak<sup>24</sup>, J. Carl Pallais<sup>18</sup>, Christina G. S. Palmer<sup>27</sup>, Jeanette C. Papp<sup>27</sup>, Neil H. Parker<sup>27</sup>, John A. Phillips III<sup>26</sup>, Jennifer E. Posey<sup>19</sup>, Lorraine Potocki<sup>19</sup>, Bradley Power<sup>16</sup>, Barbara N. Pusey<sup>16</sup>, Aaron Quinlan<sup>20</sup>, Archana N. Raja<sup>21</sup>, Deepak A. Rao<sup>18</sup>, Wendy Raskind<sup>17</sup>, Genecee Renteria<sup>27</sup>, Chloe M. Reuter<sup>21</sup>, Lynette Rives<sup>26</sup>, Amy K. Robertson<sup>26</sup>, Lance H. Rodan<sup>18</sup>, Jill A. Rosenfeld<sup>19</sup>, Natalie Rosenwasser<sup>17</sup>, Francis Rossignol<sup>16</sup>, Maura Ruzhnikov<sup>21</sup>, Ralph Sacco<sup>22</sup>, Jacinda B. Sampson<sup>21</sup>, Susan L. Samson<sup>19</sup>, Mario Saporta<sup>22</sup>, Judy Schaechter<sup>22</sup>, Timothy Schedl<sup>24</sup>, Kelly Schoch<sup>29</sup>, C. Ron Scott<sup>17</sup>, Daryl A. Scott<sup>19</sup>, Vandana Shashi<sup>29</sup>, Jimann Shin<sup>24</sup>, Rebecca H. Signer<sup>27</sup>, Edwin K. Silverman<sup>18</sup>, Janet S. Sinsheimer<sup>27</sup>, Kathy Sisco<sup>23</sup>, Edward C. Smith<sup>29</sup>, Kevin S. Smith<sup>21</sup>, Emily Solem<sup>26</sup>, Lilianna Solnica-Krezel<sup>24</sup>, Ben Solomon<sup>16</sup>, Rebecca C. Spillmann<sup>29</sup>, Joan M. Stoler<sup>18</sup>, Jennifer A. Sullivan<sup>29</sup>, Kathleen Sullivan<sup>25</sup>, Angela Sun<sup>17</sup>, Shirley Sutton<sup>21</sup>, David A. Sweetser<sup>18</sup>, Virginia Sybert<sup>17</sup>, Holly K. Tabor<sup>21</sup>, Amelia L. M. Tan<sup>32</sup>, Queenie K.-G. Tan<sup>29</sup>, Mustafa Tekin<sup>22</sup>, Fred Telischi<sup>22</sup>, Willa Thorson<sup>22</sup>, Audrey Thurm<sup>16</sup>, Cynthia J. Tiffit<sup>16</sup>, Camilo Toro<sup>16</sup>, Alyssa A. Tran<sup>19</sup>, Brianna M. Tucker<sup>21</sup>, Tiina K. Urv<sup>16</sup>, Adeline Vanderver<sup>25</sup>, Matt Velinder<sup>20</sup>, Dave Viskochil<sup>20</sup>, Tiphonie P. Vogel<sup>19</sup>, Colleen E. Wahl<sup>16</sup>, Melissa Walker<sup>18</sup>, Stephanie Wallace<sup>17</sup>, Nicole M. Walley<sup>29</sup>, Chris A. Walsh<sup>18</sup>, Jennifer Wambach<sup>23</sup>, Jijun Wan<sup>27</sup>, Lee-kai Wang<sup>27</sup>, Michael F. Wangler<sup>34</sup>, Patricia A. Ward<sup>31</sup>, Daniel Wegner<sup>23</sup>, Mark Wener<sup>17</sup>, Tara Wenger<sup>17</sup>, Katherine Wesseling Perry<sup>27</sup>, Monte Westerfield<sup>35</sup>, Matthew T. Wheeler<sup>21</sup>, Jordan Whitlock<sup>28</sup>, Lynne A. Wolfe<sup>16</sup>, Jeremy D. Woods<sup>27</sup>, Shinya Yamamoto<sup>34</sup>, John Yang<sup>16</sup>, Muhammad Yousef<sup>16</sup>, Diane B. Zastrow<sup>21</sup>, Wadiah Zein<sup>16</sup>, Chunli Zhao<sup>21</sup> and Stephan Zuchner<sup>22</sup>

<sup>16</sup>National Institutes of Health, Undiagnosed Diseases Program Clinical Site, Bethesda, MD, USA. <sup>17</sup>University of Washington and Seattle Children's Hospital Clinical Site, Seattle, WA, USA. <sup>18</sup>Harvard-affiliated Boston Children's Hospital, Massachusetts General Hospital, Brigham and Women's Hospital, and Brigham Genomic Medicine Clinical Site, Boston, MA, USA. <sup>19</sup>Baylor College of Medicine, Clinical Site, Houston, TX, USA. <sup>20</sup>University of Utah Clinical Site, Salt Lake City, UT, USA. <sup>21</sup>Stanford University Clinical Site, Stanford, CA, USA. <sup>22</sup>University of Miami Clinical Site, Miami, FL, USA. <sup>23</sup>Washington University of Saint Louis, Clinical Site, Saint Louis, MO, USA. <sup>24</sup>Washington University of Saint Louis, Model Organism Screening Center, Saint Louis, MO, USA. <sup>25</sup>Children's Hospital of Philadelphia or University of Pennsylvania Clinical Site, Philadelphia, PA, USA. <sup>26</sup>Vanderbilt University Clinical Site, Nashville, TN, USA. <sup>27</sup>University of California, Los Angeles, Clinical Site, Los Angeles, CA, USA. <sup>28</sup>University of Alabama Coordinating Center, Birmingham, AL, USA. <sup>29</sup>Duke University Clinical Site, Durham, NC, USA. <sup>30</sup>Mayo Clinic Metabolomics Core, Rochester, MN, USA. <sup>31</sup>Baylor Genetics Sequencing Core, Houston, TX, USA. <sup>32</sup>Harvard Medical School Coordinating Center, Boston, MA, USA. <sup>33</sup>Columbia University Clinical Site, New York City, NY, USA. <sup>34</sup>Baylor College of Medicine, Model Organism Screening Center, Houston, TX, USA. <sup>35</sup>University of Oregon, Model Organism Screening Center, Eugene, OR, USA.