

GENETIC DIAGNOSIS

Improving genetic diagnosis in Mendelian disease with transcriptome sequencing

Beryl B. Cummings,^{1,2,3} Jamie L. Marshall,^{1,2} Taru Tukiainen,^{1,2} Monkol Lek,^{1,2,4,5} Sandra Donkervoort,⁶ A. Reghan Foley,⁶ Veronique Bolduc,⁶ Leigh B. Waddell,^{4,5} Sarah A. Sandaradura,^{4,5} Gina L. O'Grady,^{4,5} Elicia Estrella,⁷ Hemakumar M. Reddy,⁸ Fengmei Zhao,^{1,2} Ben Weisburd,^{1,2} Konrad J. Karczewski,^{1,2} Anne H. O'Donnell-Luria,^{1,2} Daniel Birnbaum,^{1,2} Anna Sarkozy,⁹ Ying Hu,⁶ Hernan Gonorazky,¹⁰ Kristl Claeys,¹¹ Himanshu Joshi,⁵ Adam Bournazos,^{4,5} Emily C. Oates,^{4,5} Roula Ghaoui,^{4,5} Mark R. Davis,¹² Nigel G. Laing,^{12,13} Ana Topf,¹⁴ Genotype-Tissue Expression Consortium, Peter B. Kang,^{7,8} Alan H. Beggs,⁷ Kathryn N. North,¹⁵ Volker Straub,¹⁴ James J. Dowling,¹⁰ Francesco Muntoni,⁹ Nigel F. Clarke,^{4,5*} Sandra T. Cooper,^{4,5} Carsten G. Bönnemann,⁶ Daniel G. MacArthur^{1,2†}

Exome and whole-genome sequencing are becoming increasingly routine approaches in Mendelian disease diagnosis. Despite their success, the current diagnostic rate for genomic analyses across a variety of rare diseases is approximately 25 to 50%. We explore the utility of transcriptome sequencing [RNA sequencing (RNA-seq)] as a complementary diagnostic tool in a cohort of 50 patients with genetically undiagnosed rare muscle disorders. We describe an integrated approach to analyze patient muscle RNA-seq, leveraging an analysis framework focused on the detection of transcript-level changes that are unique to the patient compared to more than 180 control skeletal muscle samples. We demonstrate the power of RNA-seq to validate candidate splice-disrupting mutations and to identify splice-altering variants in both exonic and deep intronic regions, yielding an overall diagnosis rate of 35%. We also report the discovery of a highly recurrent de novo intronic mutation in *COL6A1* that results in a dominantly acting splice-gain event, disrupting the critical glycine repeat motif of the triple helical domain. We identify this pathogenic variant in a total of 27 genetically unsolved patients in an external collagen VI-like dystrophy cohort, thus explaining approximately 25% of patients clinically suggestive of having collagen VI dystrophy in whom prior genetic analysis is negative. Overall, this study represents a large systematic application of transcriptome sequencing to rare disease diagnosis and highlights its utility for the detection and interpretation of variants missed by current standard diagnostic approaches.

INTRODUCTION

The advent of whole-exome sequencing (WES) and whole-genome sequencing (WGS) has greatly accelerated our capacity to identify variants that explain many Mendelian diseases in both known and new disease genes. Although these technologies are mainstays in Mendelian

disease diagnosis, their success rate for detecting causal variants is far from complete, ranging from 25 to 50% (1–4). The primary challenge of these genome-based diagnostics is that the capacity of WES and WGS to discover genetic variants substantially exceeds our ability to interpret their functional and clinical impact (5–7).

One approach to improve the interpretation of genetic variation is to integrate functional genomic information such as RNA sequencing (RNA-seq), which provides direct insight into transcriptional perturbations caused by genetic changes (8, 9). Analysis of the complementary DNA (cDNA) of single genes has proven useful on a case-by-case basis to provide diagnoses to patients with Mendelian disorders (10–13), and RNA-seq has previously been used to observe the effect of pathogenic variants, which were identified through DNA sequencing (14, 15). However, the use of transcriptome sequencing has not yet been assessed for the discovery of pathogenic variants in a cohort of Mendelian disease patients. Such approaches have already proven useful for elucidating mechanisms of cancer and common disease (16, 17) but are not currently systematically applied to rare disease diagnosis.

Here, we describe the application of this technology to the diagnosis of patients with a range of primary muscle disorders, including myopathies and muscular dystrophies, using RNA obtained from affected muscle tissue (table S1). To investigate the value of RNA-seq for diagnosis, we obtained primary muscle RNA from 63 patients with putatively monogenic muscle disorders. Thirteen of these cases had been previously diagnosed with variants expected to have an effect on transcription, such as loss-of-function or essential splice site variants, allowing us to validate the capability of RNA-seq to identify transcriptional aberrations (table S2).

¹Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA 02114, USA. ²Medical and Population Genetics, Broad Institute of Massachusetts Institute of Technology and Harvard, Cambridge, MA 02142, USA. ³Program in Biological and Biomedical Sciences, Harvard Medical School, Boston, MA 02115, USA. ⁴School of Paediatrics and Child Health, University of Sydney, Sydney, New South Wales 2006, Australia. ⁵Institute for Neuroscience and Muscle Research, Kids Research Institute, The Children's Hospital at Westmead, Sydney, New South Wales 2145, Australia. ⁶Neuromuscular and Neurogenetic Disorders of Childhood Section, Neurogenetics Branch, National Institute of Neurological Disorders and Stroke, National Institutes of Health, Bethesda, MD 20892, USA. ⁷Division of Genetics and Genomics, Manton Center for Orphan Disease Research, Boston Children's Hospital, Harvard Medical School, Boston, MA 02115, USA. ⁸Division of Pediatric Neurology, Department of Pediatrics, University of Florida College of Medicine, Gainesville, FL 32610, USA. ⁹Dubowitz Neuromuscular Centre, University College London Institute of Child Health, London WC1N 1EH, U.K. ¹⁰Division of Neurology, Hospital for Sick Children, Toronto, Ontario M5G 1X8, Canada. ¹¹Department of Neurology, University Hospitals Leuven and University of Leuven (Katholieke Universiteit Leuven), Leuven 3000, Belgium. ¹²Department of Diagnostic Genomics, PathWest Laboratory Medicine, Perth, Western Australia 6009, Australia. ¹³Harry Perkins Institute of Medical Research, University of Western Australia, Perth, Western Australia 6009, Australia. ¹⁴John Walton Muscular Dystrophy Research Centre, MRC (Medical Research Council) Centre for Neuromuscular Diseases, Institute of Genetic Medicine, Newcastle University, Newcastle upon Tyne NE1 3BZ, U.K. ¹⁵Murdoch Childrens Research Institute, Royal Children's Hospital, Parkville, Melbourne, Victoria 3052, Australia.

*Deceased.

†Corresponding author. Email: danmac@broadinstitute.org

The remaining cohort of 50 genetically undiagnosed patients included cases for whom DNA sequencing had prioritized variants predicted to alter RNA splicing or strong candidate genes, as well as cases with no strong candidates from genetic analysis (see Fig. 1A and Materials and Methods for inclusion criteria).

RESULTS

Importance of sequencing the disease-relevant tissue

Recent large-scale studies have shown that gene expression and mRNA isoforms vary widely across tissues, indicating that for many diseases, sequencing the disease-relevant tissue will be valuable for the correct interpretation of genetic variation (18, 19). This is illustrated by the relative expression of known muscle disease genes in skeletal muscle, whole-blood, and fibroblast samples from the Genotype-Tissue Expression (GTEx) Consortium project (fig. S1) (20). A majority of the most commonly disrupted genes in muscle disease are poorly expressed in blood and fibroblasts, suggesting that RNA-seq from these easily accessible tissues may be underpowered to detect relevant transcriptional aberrations in certain genes. For these reasons, we chose to pursue RNA-seq from primary muscle tissue biopsies, which are routinely performed as part of the diagnostic evaluation of undiagnosed muscle disease patients (21, 22).

Comparison of patient RNA-seq to a muscle RNA-seq reference panel

Patient muscle samples were sequenced using the same protocol as in the GTEx project (20) and analyzed using identical pipelines to minimize technical differences, with patients sequenced at or above the same coverage as GTEx controls. From 430 skeletal muscle RNA-seq samples available through GTEx, we selected a subset of 184 samples based on RNA-seq quality metrics including RNA integrity score and ischemic time, as well as phenotypic features such as age, body mass index (BMI), and cause of death to more closely match our patient samples.

Comparison between our GTEx reference panel and patient muscle RNA-seq samples showed analogous quality metrics (table S3). Principal component analysis (PCA) of expression and splicing profiles demonstrated that patient muscle RNA-seq closely resembled control muscle when compared to tissues that potentially contaminate muscle biopsies, such as skin or fat, despite variation in the site of muscle biopsy across patients (Fig. 1B, fig. S2A, and table S1). On the basis of this clustering, we removed two samples from analysis because their expression patterns clustered more closely with GTEx adipose tissue than muscle, consistent with tissue contamination or late-stage degenerative muscle pathology (fig. S2B). We also performed fingerprinting on patient WES, WGS, and RNA-seq data to ensure that the

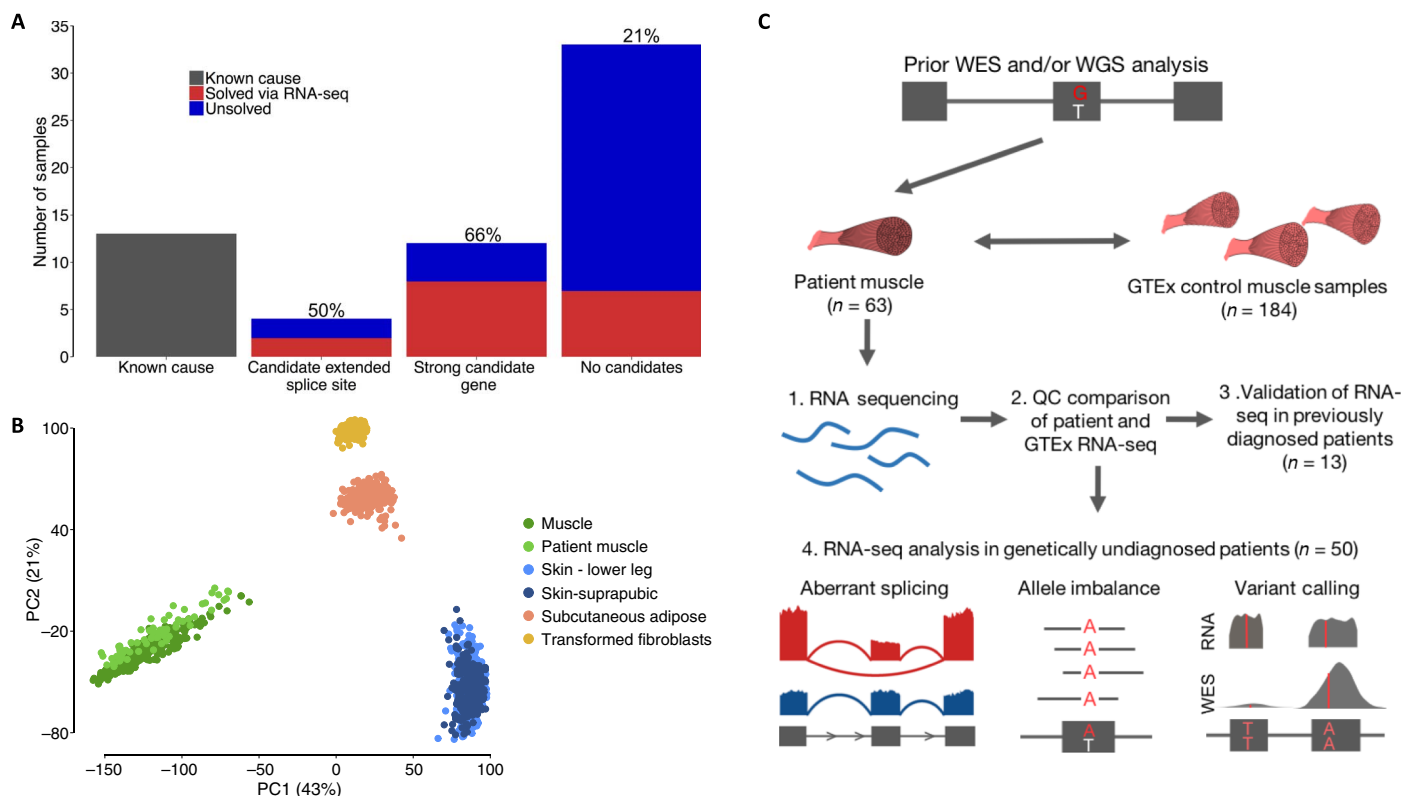


Fig. 1. Experimental design and quality control. (A) Overview of the number of samples that underwent RNA-seq. We performed RNA-seq on 13 previously genetically diagnosed patients, 4 patients in whom previous genetic analysis had identified an extended splice site variant of unknown significance (VUS), 12 patients in whom genetic analysis had identified a strong candidate gene, and 34 patients with no strong candidates from previous analysis. RNA-seq enabled the diagnosis of 35% of patients overall, with the rate, shown above the bar plots, varying depending on previous evidence from genetic analysis. (B) PCA based on gene expression profiles of patient muscle samples passing quality control (n = 61) and GTEx samples of tissues that potentially contaminate muscle biopsies shows that patient samples cluster closely with GTEx skeletal muscle. (C) Overview of experimental setup and RNA-seq analyses performed. Our framework is based on identifying transcriptional aberrations that are present in patients and missing in GTEx controls. Upon ensuring that GTEx and patient RNA-seq data were comparable, we validated the capacity of RNA-seq to resolve transcriptional aberrations in previously diagnosed patients and performed analyses of aberrant splicing, allele imbalance, and variant calling in our remaining cohort of genetically undiagnosed muscle disease patients.

source of DNA sequencing and muscle RNA-seq data was the same individual.

We explored the utility of analyzing patient RNA-seq data to detect aberrant splice events and allele-specific expression and performed variant calling from RNA-seq data to identify pathogenic events or to prioritize genes for closer analysis (Fig. 1C). We also identified outlier gene expression status in patients; however, this analysis was underpowered to prioritize candidate genes in our study (fig. S3). The resulting diagnoses were made primarily through the detection of aberrant splice events in patients, with information on gene-level allele imbalance playing a complementary role.

In previously diagnosed cases, manual evaluation of pathogenic essential splice site variants revealed a splice aberration, such as exon skipping or extension, demonstrating that RNA-seq can help resolve the effect of variants on transcription (fig. S4, A to F). To detect aberrant transcriptional events genome-wide, we developed an approach based on identifying high-quality exon-exon splice junctions present in patients or groups of patients and missing in GTEx controls (code available at <https://github.com/berylc/MendelianRNA-seq>). We performed splice junction discovery from split-mapped reads, considering only those that were uniquely aligned and nonduplicate. To account for library size and stochastic gene expression differences between samples, we performed local normalization of read counts based on read support for overlapping annotated junctions (fig. S5, A and B). We then performed filtering of splice junctions based on the number of samples in which a splice junction is observed and the number of reads and normalized value supporting that junction in each sample. Our approach successfully reidentified all known pathogenic events in patients in whom manual evaluation had revealed aberrant splicing around splice variants previously identified through genomic testing. We defined filtering parameters that selectively identified these previously known aberrant splice events and applied them to our remaining cohort of undiagnosed patients. This method resulted in the identification of a median of 5, 26, and 190 potentially pathogenic splice events per sample in ~190 neuromuscular disease associated genes, Online Mendelian Inheritance in Man (OMIM) genes, and all genes, respectively (fig. S6), which required manual curation to interpret pathogenicity and led to the diagnoses made in this study.

Diagnoses made via RNA-seq

RNA-seq allowed the diagnosis of 17 previously unsolved families, yielding an overall diagnosis rate of 35% in this challenging subset of rare disease patients for whom extensive prior analysis of DNA sequencing data had failed to return a genetic diagnosis. We also identified splice disruption in other known and putatively novel disease genes in several patients; however, due to unavailability of additional information, such as parental DNA, we could not pursue these cases further (fig. S7). Detection of aberrant splicing led to the identification of a broad class of both coding and noncoding pathogenic variants, resulting in a range of splice defects such as exon skipping, exon extension, and exonic and intronic splice gain, which were validated by reverse transcription polymerase chain reaction (RT-PCR) analysis (see Fig. 2, Table 1, and the Supplementary Materials and Methods). RNA-seq patterns also helped pinpoint three structural variants in *DMD* that were subsequently confirmed by WGS (fig. S8).

Cases diagnosed in this study highlight several key advantages of RNA-seq in rare disease diagnosis to confirm the pathogenicity of variants and to detect previously unidentified variation. In four patients

with previously detected extended splice site VUS, RNA-seq confirmed splice disruption in two patients (Fig. 1A and fig. S9, A and B). The variants had no observable effect on local splicing patterns in the remaining two patients, emphasizing the value of RNA-seq in ruling out non-pathogenic VUS (fig. S9, C and D).

RNA-seq also led to the identification of an additional disruptive extended splice site variant missed by exome sequencing. In a nemaline myopathy patient with one previously detected recessive frameshift variant in the *NEB* gene, RNA-seq identified an exon extension event caused by an underlying variant at the +3 position of the donor site, which led to the introduction of a premature stop codon to the transcript as the second recessive allele (Fig. 2B). The exon harboring this variant was not captured in the exome kit used to screen the patient (fig. S10), underlining the utility of RNA-seq at complementing WES to identify previously undetected variants.

Synonymous and missense variants in large, variation-rich genes, such as *TTN*, are exceptionally challenging to interpret and are often filtered out in DNA sequencing pipelines (23, 24). With RNA-seq, we were able to assign pathogenicity to a missense variant in *TTN* and two synonymous variants in *RYR1* and *POMGNT1* (fig. S11). In patient N22, the identified missense variant created a GT donor splice site for which the consensus motif included a G nucleotide in the +5 position, known to contribute to the strength of the splice site (25, 26). The well-conserved donor +5-G motif was missing in the competing canonical splice site, thus resulting in a stronger novel splice site and gain of splicing from the exon body (Fig. 2C). A similar mechanism was observed in *RYR1*, caused by a synonymous variant in a patient carrying a second pathogenic allele in the gene (fig. S11A). In an additional patient carrying an essential splice site variant in *POMGNT1*, we identified a synonymous variant disrupting an exonic splice motif and resulting in exon skipping (fig. S11, B to D).

In eight cases, RNA-seq aided in the identification of noncoding pathogenic variants. We identified splice site–creating hemizygous deep intronic variants in *DMD* that resulted in the creation of a pseudoexon and led to a premature stop codon in the coding sequence in three patients (Fig. 2D and fig. S12). Although RNA-seq from a patient with severe Duchenne muscular dystrophy showed only splicing to the pseudoexon (fig. S12), wild-type splicing between annotated exons was observed in two patients with a milder Becker muscular dystrophy phenotype, indicating the presence of residual functional *DMD* transcripts that explain the milder disease course. Such intronic variants are unobservable with WES and too abundant to be interpretable with WGS alone, emphasizing the utility of RNA-seq at resolving pathogenicity of these noncoding variants.

In two patients with no strong candidates from WES and WGS (N22 and N25), we identified heterozygous splice disruption in two commonly disrupted recessive muscle disease genes, *NEB* and *TTN*. These genes harbor regions with highly similar sequences, the so-called triplicate repeat regions (27, 28). Because of high sequence similarity, the region has poor mapping quality, resulting in low-quality variant calls that are filtered by the most current diagnostic pipelines. To identify possible pathogenic variants in the triplicated regions of *NEB* and *TTN* in these two patients, we developed a method based on remapping the triplicate regions to a detriplicated pseudoreference and performing hexaploid variant calling (fig. S13, A to C). This method was applied to available WES/WGS and RNA-seq data for all patients and identified one novel nonsense and one novel frameshift variant in *NEB* and *TTN* in these two patients, which finalized their diagnoses (fig. S13D, N25, and fig. S13E, N22).

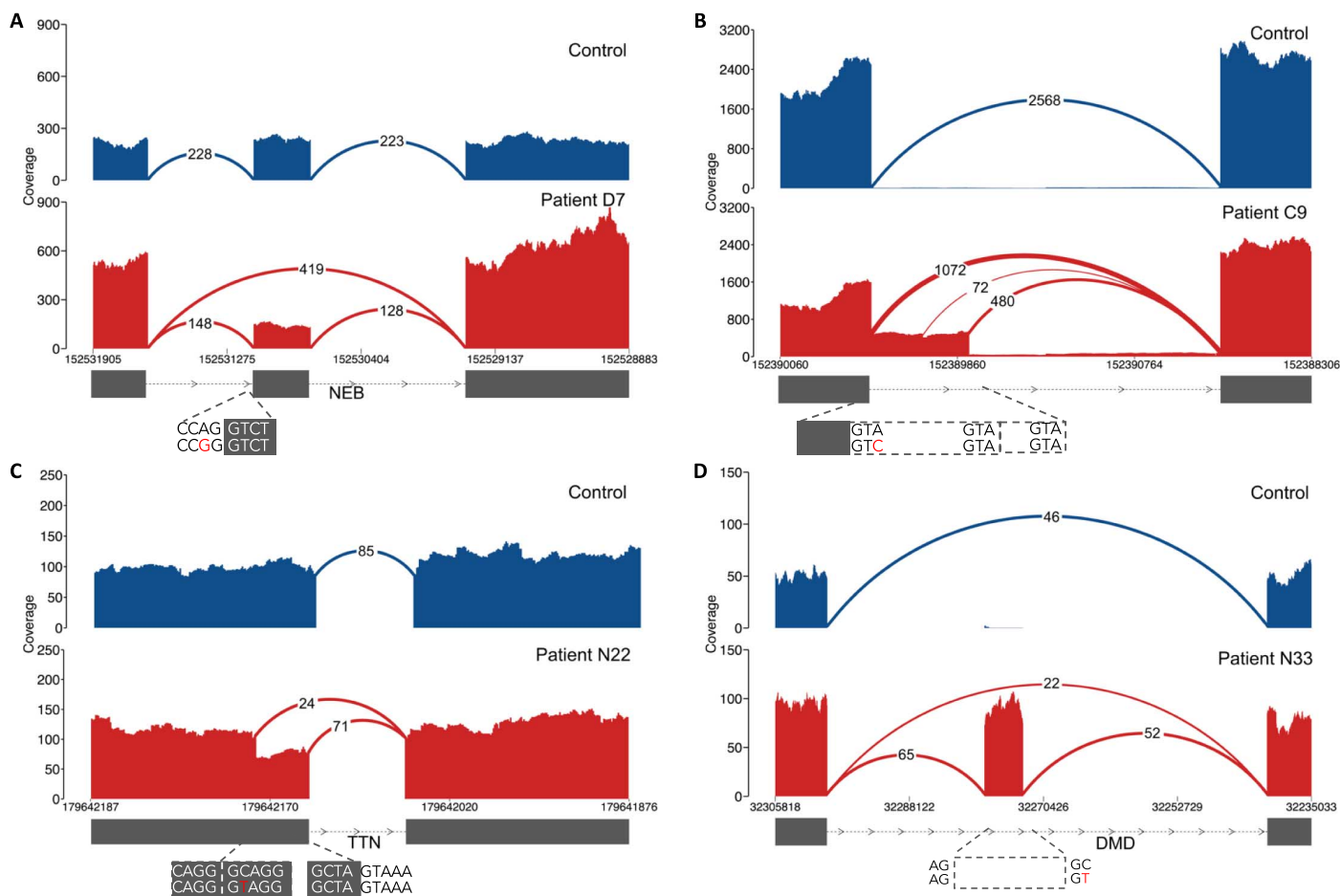


Fig. 2. Types of pathogenic splice aberrations discovered in patients. RNA-seq identified a range of aberrations caused by both coding and noncoding variants, such as (A) exon skipping caused by an essential splice site variant in patient D7, (B) exon extension caused by a donor +3 A>C extended splice site variant in nemaline myopathy patient C9 (where disruption of splicing at the canonical splice site results in splicing from intact GTA motifs from the intron), (C) exonic splice gain caused by a C>T donor splice site–creating variant in patient N22 with a donor +5-G sequence context, resulting in a stronger splice motif than the existing canonical splice site, and (D) intronic splice gain in patient N33 caused by a C>T donor splice site–creating deep intronic variant. Evidence for wild-type splicing in addition to the inclusion of the pseudoexon in the patient is in line with the milder Becker's muscular dystrophy phenotype. Splice aberrations shown in (B) to (D) result in the introduction of a premature stop codon to the transcript.

Identification of a recurrent splice site–creating variant in collagen VI–related dystrophy

A notable example of the power of transcriptome sequencing is our discovery of a genetic subtype of severe collagen VI–related dystrophy, which is caused by mutations in one of the three collagen VI genes (*COL6A1*, *COL6A2*, and *COL6A3*) (21). In four patients who had previously tested negative with deletion/duplication testing and fibroblast cDNA sequencing of the collagen VI genes as well as clinical WES and WGS, we identified an intron inclusion event in *COL6A1* using RNA-seq (Fig. 3A). The splicing-in of this intronic segment, which is missing in GTEx controls and all other patients in our cohort, is caused by a donor splice site–creating GC>GT variant that pairs with a cryptic acceptor splice site 72 base pairs (bp) upstream, creating an in-frame pseudoexon (Fig. 3B). This variant is missing in the 1000 Genomes Project data set (29) as well as an in-house data set of 5500 control WGS samples. The resulting inclusion of 24 amino acids occurs within the N-terminal triple-helical collagenous G-X-Y repeat region of the *COL6A1* gene, the disruption of which has been well established to cause dominant-negative pathogenicity in a variety of collagen disor-

ders (30). Notably, cDNA analysis shows that the aberrant transcript is observable in muscle but in much smaller amounts in cultured dermal fibroblasts, making the event identifiable by muscle transcriptome analysis despite being previously missed by fibroblast cDNA sequencing (Fig. 3C). Using this information, we genotyped the variant in a larger, genetically undiagnosed collagen VI–like dystrophy cohort and identified 27 additional patients carrying the intronic variant. We confirmed that the variant had occurred as an independent de novo mutation in all 16 families for whom trio DNA was available. On the basis of this screening, we estimate that up to a quarter of all cases clinically suggestive of collagen VI–related dystrophy but negative by exon-based sequencing are due to this recurrent de novo mutation (see the Supplementary Materials and Methods).

Evaluation of splice prediction algorithms and RNA-seq in alternative tissues

Exons harboring the pathogenic variants identified in this study show low coverage in GTEx whole-blood and fibroblast samples, indicating that a majority of these diagnoses likely could not have been made using

Table 1. Diagnoses made in the study via patient muscle RNA-seq.					
Patient	Phenotype	Gene	Variants	Variant class	Effect
E2	Nemaline myopathy	NEB	chr2: 152,544,805 C>T chr2: 152,520,057 C>T	Essential splice, extended splice	Exon skipping + exon extension, exon extension
C9	Nemaline myopathy	NEB	chr2: 152,581,432 TG>T chr2: 152,389,953 A>C	Frameshift, extended splice	Exon extension
E4	Fetal akinesia	TTN	chr2: 179,586,600 CAT>C chr2: 179,446,219 ATACT>A	Frameshift, extended splice	Exon skipping
C6	Duchenne muscular dystrophy	DMD	chrX: 32,366,860 A>C	Intronic variant	Intronic splice gain
N33	Myalgia, myoglobinuria	DMD	chrX: 32,274,692 G>A	Intronic variant	Intronic splice gain
C7	Becker muscular dystrophy	DMD	chrX: 31,613,687 G>T	Intronic variant	Intronic splice gain
N29	Collagen VI-related dystrophy	COL6A1	chr21: 47,409,881 C>T	Intronic variant	Intronic splice gain
N30	Collagen VI-related dystrophy	COL6A1	chr21: 47,409,881 C>T	Intronic variant	Intronic splice gain
N31	Collagen VI-related dystrophy	COL6A1	chr21: 47,409,881 C>T	Intronic variant	Intronic splice gain
N32	Collagen VI-related dystrophy	COL6A1	chr21: 47,409,881 C>T	Intronic variant	Intronic splice gain
N25	Nemaline myopathy	NEB	chr2: 152,355,017 G>T chr2: 152,449,646G>A	Intronic variant, nonsense	Intronic splice gain
C11	Congenital fiber-type disproportion	RYR1	chr19: 38,958,362 C>T chr19: 38,958,372 G>A	Synonymous, missense	Exonic splice gain
N22	Multi/minicore congenital myopathy	TTN	chr2: 179,642,185 G>A chr2: 179,523,240 CTCT>C	Missense, frameshift	Exonic splice gain
C1	α -Dystroglycanopathy	POMGNT1	chr1: 46,655,129 C>A chr1: 46,660,532 G>A	Essential splice, synonymous	Exonic splice gain, exon skipping
C3	Duchenne muscular dystrophy	DMD	chrX: 31,790,694–31,798,498	Inversion-deletion	Exon skipping
C2	Duchenne muscular dystrophy	DMD	chrX: 31,378,946–151,194,962	Inversion	Splice disruption
C4	Duchenne muscular dystrophy	DMD	chrX: 32,521,820–35,180,380	Inversion	Splice disruption

RNA-seq from these tissues (fig. S14). Furthermore, many of the diagnoses made in this study could not have been made on genotype information alone, because splice prediction algorithms alone are currently insufficient to classify variants as causal (31, 32). Although existing in silico algorithms correctly predicted disruption for the two extended splice site VUS in our study, they also generated false-positive predictions for the remaining two extended splice site variants with no effect on splicing (see fig. S15A and the Supplementary Materials and Methods). In addition, existing algorithms showed poor specificity in identifying splice site–creating coding variants, identifying on average more than 100 putative splice site–creating rare variants [$<1\%$ population frequency in Exome Aggregation Consortium (ExAC)] exome-wide (fig. S15B).

DISCUSSION

Our results show that RNA-seq is valuable for the interpretation of coding as well as noncoding variants and can provide a substantial increase in diagnosis rate in patients for whom exome or whole-genome analysis has not yielded a molecular diagnosis. In our cohort, RNA-seq led to the diagnosis of 66% of patients where clinical phenotyping and

DNA sequencing prioritized a strong candidate gene. In comparison, through identifying aberrant splice events found in patients and missing in GTEx controls, we were able to diagnose 21% of patients with no strong candidates from WGS or WES.

Our work illustrates the value of large multitissue transcriptome data sets such as GTEx to serve as a reference to facilitate the identification of extreme splicing or allele balance outlier events in patients. In the case of muscle disorders, our diagnoses were made primarily through direct identification of aberrations in splicing using the GTEx skeletal muscle RNA-seq data set as a reference panel. Our present work focused on identifying such aberrations in known muscle disease genes, and the considerably lower number of putatively pathogenic events identified in neuromuscular disease genes versus all genes underlines the advantage of a candidate gene list for this analysis. Further improvements in filtering identified splice junctions to obtain a smaller list of candidate events will be useful to expand this work for new disease gene discovery. In addition, with increasing sample sizes and improvements in methods, RNA-seq can also be used to identify somatic variants and to detect regulatory variants upstream, through analysis of expression status and allelic imbalance.

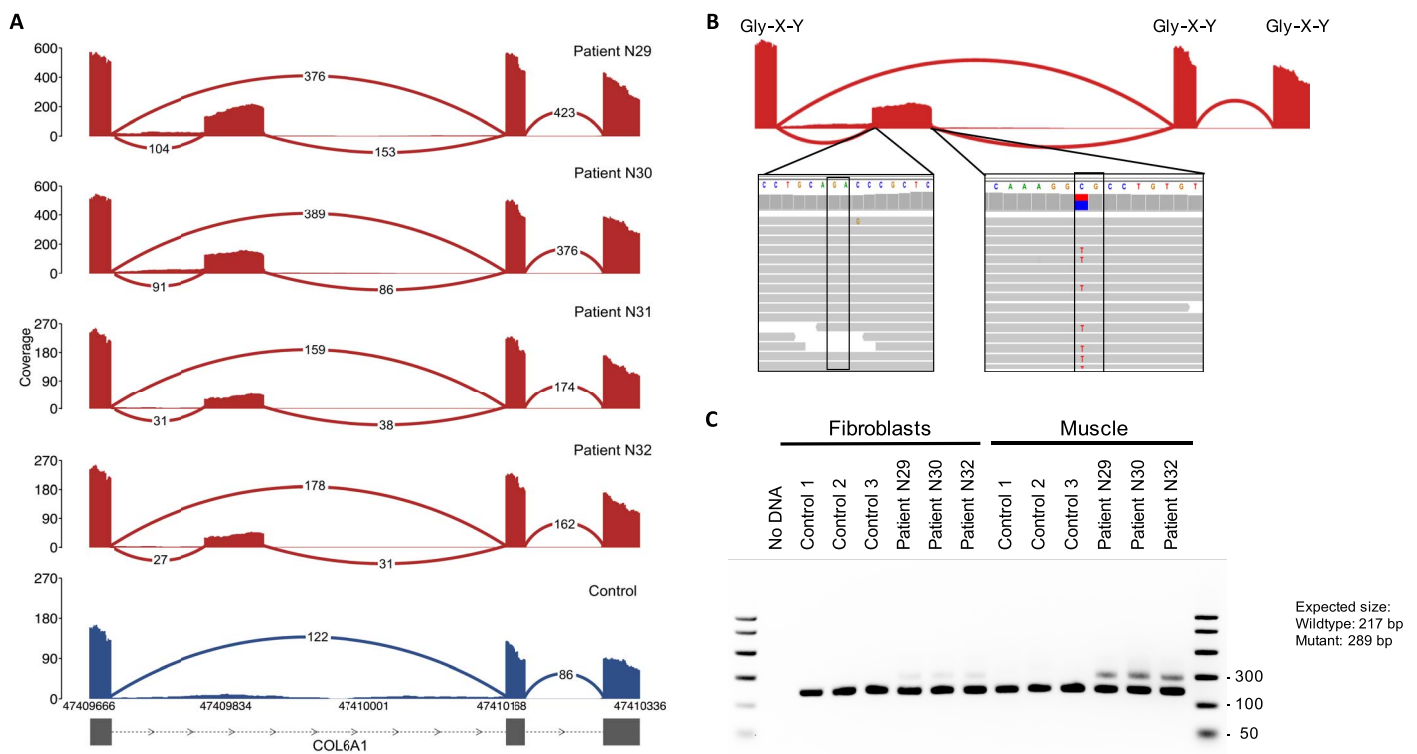


Fig. 3. Identification of a recurrent splice site-creating variant in four collagen VI-related dystrophy patients. (A) Splicing-in of the pseudoexon was observed in four patients in our cohort (red) and missing in all other patients and GTEx samples (blue). (B) Inclusion of the 24-amino acid segment is caused by a C>T donor splice site-creating variant, which pairs with an AG splice acceptor site 72 bp upstream. The variant is found in a CpG nucleotide context, which likely explains its recurrent de novo status, and disrupts the Gly-X-Y repeat motifs of COL6A1. (C) The inclusion event is observable in RT-PCR amplicons from patient muscle but is found at comparatively lower levels in cultured dermal fibroblasts derived from the patients, explaining why the pathogenic event was missed in all four patients through previous fibroblast cDNA sequencing.

Access to the disease-relevant tissue for many Mendelian disorders remains a major barrier for the use of transcriptome sequencing in genetic diagnosis. The RNA-seq framework developed in this study can be adapted for rare diseases where biopsies are available, such as Mendelian disorders affecting the heart, kidney, liver, skin, and other tissues. For example, during the preparation of this paper, the application of RNA-seq to fibroblast samples for the genetic diagnosis of mitochondrial disease was reported in an unpublished preprint (33). For disorders where biopsy of the disease-relevant tissue is unattainable, analyses are possible through identification of proxy tissues using databases such as GTEx and careful consideration of the expression status of the relevant genes in the proxy tissue. Alternatively, the framework developed in this study can also enable diagnoses through reprogramming patient cells into induced pluripotent stem cells and differentiation into disease-relevant tissues of interest.

Evaluation of existing splice prediction algorithms for the splice-disrupting variants identified in the study highlights that information on DNA sequence alone does not currently match the ability of RNA-seq to identify the transcriptional consequences of variants on a genome-wide scale. The diagnoses made in our study with RNA-seq, particularly the discovery of the highly recurrent mutation in COL6A1, demonstrate that other such cryptic splice-affecting variants may contribute substantially to undiagnosed diseases that have evaded prior detection with exome or whole-genome analysis. Overall, this work suggests that RNA-seq is a valuable component of the diagnostic toolkit for rare diseases and can aid in the identification of new pathogenic variants in known genes as well as new mechanisms for Mendelian disease.

MATERIALS AND METHODS

Study design

We sought to explore the utility of transcriptome sequencing as a complementary diagnostic tool to exome and whole-genome analysis. We reasoned that RNA-seq would allow us to interpret variants previously identified through genetic analysis and may pinpoint genetic lesions that may have eluded DNA sequencing. To interpret transcriptional aberrations seen in patients, we obtained a reference panel of 184 sets of skeletal muscle RNA-seq data from the GTEx project. Our framework was based on identifying transcriptional aberrations present in patients but missing in GTEx controls. We first validated the capacity of RNA-seq to resolve transcriptional aberrations in 13 patients with prior genetic diagnosis and then analyzed the remaining 50 genetically undiagnosed patients to detect aberrant splice events and allele-specific expression and performed variant calling from RNA-seq data to identify pathogenic events or to prioritize genes for closer analysis.

Clinical sample selection

Patient cases with available muscle biopsies were referred by clinicians from March 2013 through June 2016. Samples fell into four broad categories:

(1) Patients for whom previous genetic analysis had resulted in a diagnosis with at least one loss-of-function or essential splice site variant, serving as positive controls to assess the capability of RNA-seq to identify the transcriptional effect of the variants ($n = 13$; patient IDs starting with “D”).

(2) Patients with candidate extended splice site variants that had been categorized as VUS, for which assignment of pathogenicity would result in a complete diagnosis for the patient ($n = 4$; patient IDs starting with “E”).

(3) Patients for whom a strong candidate gene was implicated because of either a well-defined monogenic disease phenotype, such as patients with clear Duchenne muscular dystrophy evidenced by clinical diagnosis and loss of dystrophin expression ($n = 6$), or the presence of one pathogenic heterozygous variant identified in a gene matching the patient’s phenotype, without a second pathogenic variant in that gene ($n = 6$; patient IDs starting with “C”).

(4) Patients with no strong candidates based on previous genetic analysis such as WES or WGS ($n = 34$; patient IDs starting with “N”).

Patients who fit categories 2 to 4 are referred to as undiagnosed before RNA-seq and constitute the denominator for the 35% diagnosis rate. All patients had prior analysis of WES and/or WGS data, except two cases (patients E4 and D11) for whom targeted sequencing had identified candidate extended and essential splice site variants, respectively. We favored cases with previous trio WES or WGS: 29 of 63 patients had complete trios, with 3 additional patients having one parent sequenced. Although age of onset was not considered as an exclusion criterion, most of the patients in the cohort had a congenital or early childhood-onset primary muscle disorder.

Muscle biopsies or RNA were shipped frozen from clinical centers via a liquid nitrogen dry shipper and stored in liquid nitrogen cryogenic storage. Before submission to the sequencing platform, all muscle samples were visually inspected, photographed, cut into 50- μ m sections on a Leica CM1950 model cryostat, and transferred to prechilled cryotubes in preparation for RNA extraction. When muscle arrived embedded in optimum cutting temperature compound, 8- μ m transverse cryosections were mounted on positively charged Superfrost Plus slides (VWR, 48311–703) and stained with hematoxylin and eosin (H&E) to assess the relative proportion of muscle versus fibrosis and adipose infiltration as well as the presence of overt freeze-thaw artifact. All samples analyzed with H&E showed muscle quality sufficient to proceed to RNA-seq.

RNA sequencing

RNA was extracted from muscle biopsies via the miRNeasy Mini Kit from QIAGEN according to the kit’s instructions. All RNA samples were measured for quantity and quality. Samples had to meet the minimum cutoff of 250 ng of RNA and RNA quality score (RQS) of 6 to proceed with RNA-seq library preparation. A fraction of samples falling below an RQS of 6 were also submitted for sequencing. All samples submitted had a range of RQs between 3.5 and 8.

Sequencing was performed at the Broad Institute Genomics Platform using the same non-strand-specific protocol with poly-A selection of mRNA (Illumina TruSeq) used in the GTEx sequencing project (20) to ensure consistency of our samples with GTEx control data. Paired-end 76-bp sequencing was performed on Illumina HiSeq 2000 instruments, with sequence coverage of 50 million or 100 million reads. One sample (patient N33) was sequenced to a higher depth at 500 million reads to permit downsampling analysis of the effects of increasing RNA-seq depth.

Selection of GTEx controls

GTEx data were downloaded from the Database of Genotypes and Phenotypes (dbGaP) (www.ncbi.nlm.nih.gov/gap) under accession phs000424.v6.p1. From 430 available GTEx skeletal muscle RNA-seq

samples, we selected 184 samples on the basis of RNA integrity score (between 6 and 9), number of nonduplicate uniquely mapped read pairs (between 35 million and 75 million reads), and ischemic time (<12 hours) to remove any samples that were outliers for these quality metrics. GTEx samples were further filtered to remove those with known clinical conditions such as Klinefelter’s syndrome or those for whom death followed after long- or intermediate-term illness or medical intervention (Hardy scale 0, 3, or 4). Overall, approximately 80% of GTEx samples with available muscle RNA-seq are older than 40 (median age, 54) and have a BMI over 25 (median BMI, 27). Thus, we selected samples to enrich for younger GTEx donors to more closely match our patient cohort. All samples younger than 50 were selected, resulting in 76 samples with high-quality RNA-seq data. We then added older samples back on the criterion that their BMI was below 30. This resulted in a total of 184 GTEx control samples for our reference panel, with comparable male and female sample count (105 males and 79 females). This filtering method also enriched the RNA-seq data from organ donors and surgical donors as opposed to postmortem samples (72% of selected GTEx controls are derived from surgical or organ donors versus 45% in the unfiltered data set). A full list of GTEx sample IDs used as the reference panel can be found in table S4.

RNA-seq alignment and quality control

GTEx BAM files downloaded from dbGaP were realigned after conversion to FASTQ files with Picard SamToFastq. Both patient and GTEx reads were aligned via STAR 2-Pass version v.2.4.2a using hg19 as the genome reference and GENCODE V19 annotations. Briefly, first-pass alignment was performed for novel junction discovery, and the identified junctions were filtered to exclude unannotated junctions with less than five uniquely mapped read supports, as well as junctions found on the mitochondrial genome. These junctions were then used to create a new annotation file, and second-pass alignment was performed as recommended by the STAR manual to enable sensitive junction discovery. Duplicate reads were marked with Picard MarkDuplicates (v.1.1099).

Quality metrics for patient and GTEx RNA-seq data were obtained by running RNA-SeQC (v1.1.8) on STAR-aligned BAM files (34). PCA on gene expression was performed on the basis of RPKM (reads per kilobase of transcript per million mapped reads) values calculated by RNA-SeQC. Two samples (D6 and N3) were removed because of outlier status in PCA, consistent with a high proportion of nonmuscle tissue in the samples (fig. S2B). For GTEx samples, the expression and exon-level read count data were downloaded from dbGaP under accession phs000424.v6. For PCA of exon inclusion metrics, we obtained PSI (percentage spliced in) values for GTEx samples as described in (35).

To ensure that patient DNA and RNA data were identity-matched, we compared variants identified in WES, WGS, and RNA-seq data. WES, WGS, and RNA-seq data were joint-genotyped for a set of ~5800 common single nucleotide polymorphisms (SNPs) collated by Purcell *et al.* (36) using the Genome Analysis Toolkit (GATK) HaplotypeCaller package version 3.4. We then calculated pairwise inheritance by descent estimates between DNA sequencing and RNA-seq data using PLINK (v1.08p). Relatedness coefficients for WES, WGS, and RNA-seq data from the same individual ranged from 0.67 to 1.00 across our samples (mean, 0.9), compared to a range of 0 to 0.18 (mean, 0.001) for non-matching individuals, confirming that the sources for DNA sequencing and RNA-seq were the same for each patient in our data set.

Exome sequencing and WGS

WES on DNA samples (>250 ng of DNA, at >2 ng/ μ l) was performed using Illumina or Agilent SureSelect v2 exome capture. The exome

sequencing pipeline included sample plating, library preparation (2-plexing of samples per hybridization), hybrid capture, sequencing (76-bp paired reads), and sample identification quality control check. Hybrid selection libraries covered >80% of targets at 20× with a mean target coverage of >80×. The exome sequencing data were demultiplexed, and each sample's sequence data were aggregated into a single Picard BAM file. WGS was performed on 500 ng to 1.5 µg of genomic DNA using a PCR-free protocol. These libraries were sequenced on the Illumina HiSeq X10 with 151-bp paired-end reads and a target mean coverage of >300×.

Exome and genome sequencing data were processed through a Picard-based pipeline using base quality score recalibration and local realignment at known insertions/deletions (indels). The Burrows-Wheeler Aligner was used for mapping reads to the human genome build 37 (hg19). SNPs and indels were jointly called across all samples using GATK HaplotypeCaller. Default filters were applied to SNP and indel calls using the GATK variant quality score recalibration, and variants were annotated using Variant Effect Predictor (v78); additional information on this pipeline is provided in the first supplementary section of (37). The variant call set was uploaded to the seqr analysis platform (seqr.broadinstitute.org) to perform variant filtering using inheritance patterns, functional annotation, and variant frequency in reference databases including ExAC (37) and 1000 Genomes (29).

Identification of pathogenic splice events

Splice junctions were identified from split-mapped reads, considering only uniquely aligned, nonduplicate reads that passed platform/vendor quality controls. For each splice junction, we noted the following:

- (1) the genomic coordinates
- (2) the gene in which the junction was observed based on GENCODE v.19
- (3) the number of samples in which the splice junction was observed
- (4) the number of total reads supporting the junction in 245 samples (184 GTEx and 61 patient samples)
- (5) the per-sample read support for the junction.

We then performed local normalization of per-sample read support on the basis of the support for the highest shared annotated junction (fig. S5A). For example, an exon-skipping event harbors two annotated exon-intron junctions, and we normalized this by the maximum of read count support for canonical splicing at these two wild-type junctions. This local normalization allows for filtering low-level mapping noise and accounts for stochastic gene expression and library size differences between samples (fig. S5B).

To identify pathogenic splice events, splice junctions in protein coding genes were filtered in terms of the number of samples a splice junction is present in and the number of reads and the normalized value supporting that junction. Specifically, we defined a sensitive cutoff at which an aberrant splice event is seen with at least 5% of the read support as compared to the shared annotated junction, with at least two reads supporting the event. We also required a splice junction to contain at least one annotated exon-exon junction, indicating that the event was spliced into an existing transcript (fig. S5A). We performed analysis on a per-sample basis, each time requiring the normalized value of a given splice junction to be maximum in that sample and twice that of the next highest sample, allowing us to search for unique events in the patient.

All candidate pathogenic splice events were manually evaluated using the Integrative Genomics Viewer. This resulted in the identification of aberrant splicing at eight of nine pathogenic essential splice site variants and resulted in the diagnosis of 10 of 17 patients in the study. A splice

aberration was not observed around an essential splice site variant found in *TTN* in patient D5 because of insufficient number of reads mapping to the local region (fig. S4E). We extended filtering parameters to identify splice junctions present in fewer than 10 samples, but with high read support in each sample, allowing us to identify the intronic splice-gain event present in four patients in *COL6A1* (Fig. 3A). We note that this approach would also identify putatively pathogenic splice aberrations, for which there are GTEx carriers. The remaining three Duchenne muscular dystrophy patients were diagnosed through manual analysis of splicing patterns in *DMD* and resulted in the identification of splice disruption. Overlapping structural variants at these regions were confirmed by subsequent WGS (fig. S8).

Statistical analysis and code availability

Our approach to evaluating outlier status for allele imbalance in patients involved defining the 95% confidence interval (means \pm 2 SD) of mean allele balance in GTEx individuals for each gene and identifying patients for whom the gene-level allele balance fell outside of the range. Comparison between GTEx and patient RNA-seq data quality metrics relied on a *t* test for significance. Data processing, analysis, and figure generation were performed using scripts written in Python 2.7 and R 3.2; code for identifying and filtering splice junctions and for variant calling in the triplicate regions of *NEB* and *TTN* is available at <https://github.com/berylc/MendelianRNA-seq>.

SUPPLEMENTARY MATERIALS

www.sciencetranslationalmedicine.org/cgi/content/full/9/386/eaal5209/DC1
Materials and Methods

- Fig. S1. Expression of commonly disrupted muscle disease genes in muscle, blood, and fibroblasts.
Fig. S2. PCA based on PSI metrics and gene expression of GTEx and patient samples.
Fig. S3. Overview of results from expression outlier analysis.
Fig. S4. Evaluation of RNA-seq around pathogenic essential splice site variants previously identified by genetic analysis.
Fig. S5. Overview of splice junction filtering approach.
Fig. S6. Number of potentially pathogenic splice events identified per patient.
Fig. S7. Examples of splice disruption in patients with no diagnosis at the completion of the study.
Fig. S8. Identification of aberrant splicing overlapping structural variants with RNA-seq.
Fig. S9. Resolving the effect of extended splice site variants with RNA-seq.
Fig. S10. Coverage of exon harboring splice-disrupting variant identified in patient C9 in RNA-seq and WES.
Fig. S11. Assignment of pathogenicity to missense and synonymous variants with RNA-seq.
Fig. S12. Identification of pathogenic noncoding variants with RNA-seq.
Fig. S13. Overview of triplicate region remapping.
Fig. S14. Comparison of the number of reads aligning to exons harboring pathogenic variants identified in the study in GTEx muscle, whole blood, and fibroblast tissues.
Fig. S15. Evaluation of splice prediction algorithms.
Fig. S16. Identification of allele imbalance with RNA-seq.
Table S1. Overview of clinical cases that underwent RNA-seq.
Table S2. Summary of patients previously diagnosed by genetic analysis with variants expected to result in transcriptional aberrations and the corresponding effect seen in the RNA-seq data.
Table S3. Comparison of quality metrics between patient and GTEx RNA-seq samples showing correspondence between patients and controls.
Table S4. List of 184 GTEx control skeletal muscle RNA-seq samples.
Table S5. PCR conditions and primers used for RT-PCR validation of splice aberrations identified via RNA-seq and Sanger sequencing of cDNA.
Table S6. PCR conditions and primers used for genomic Sanger sequence validation of variants identified in patients.
References (38–46)

REFERENCES AND NOTES

1. A. Ankala, C. da Silva, F. Gualandi, A. Ferlini, L. J. Bean, C. Collins, A. K. Tanner, M. R. Hegde, A comprehensive genomic approach for neuromuscular diseases gives a high diagnostic yield. *Ann. Neurol.* **77**, 206–214 (2015).

2. Y. Yang, D. M. Muzny, F. Xia, Z. Niu, R. Person, Y. Ding, P. Ward, A. Braxton, M. Wang, C. Buhay, N. Veeraraghavan, A. Hawes, T. Chiang, M. Leduc, J. Beuten, J. Zhang, W. He, J. Scull, A. Willis, M. Landsverk, W. J. Craigen, M. R. Bekheirnia, A. Stray-Pedersen, P. Liu, S. Wen, W. Alcaraz, H. Cui, M. Walkiewicz, J. Reid, M. Bainbridge, A. Patel, E. Boerwinkle, A. L. Beaudet, J. R. Lupski, S. E. Plon, R. A. Gibbs, C. M. Eng, Molecular findings among patients referred for clinical whole-exome sequencing. *JAMA* **312**, 1870–1879 (2014).
3. J. C. Taylor, H. C. Martin, S. Lise, J. Broxholme, J.-B. Cazier, A. Rimmer, A. Kanapin, G. Lunter, S. Fiddy, C. Allan, A. R. Aricescu, M. Attar, C. Babbs, J. Becq, D. Beeson, C. Bento, P. Bignell, E. Blair, V. J. Buckle, K. Bull, O. Cais, H. Cario, H. Chapel, R. R. Copley, R. Cornell, J. Craft, K. Dahan, E. E. Davenport, C. Dendrou, O. Devuyst, A. L. Fenwick, J. Flint, L. Fugger, R. D. Gilbert, A. Goriely, A. Green, I. H. Greger, R. Grocock, A. V. Gruszczyn, R. Hastings, E. Hatton, D. Higgs, A. Hill, C. Holmes, M. Howard, L. Hughes, P. Humburg, D. Johnson, F. Karpe, Z. Kingsbury, U. Kini, J. C. Knight, J. Krohn, S. Lample, C. Langman, L. Lonie, J. Luck, D. McCarthy, S. J. McGowan, M. F. McMullin, K. A. Miller, L. Murray, A. H. Németh, M. N. Andrew, D. Nutt, E. Ormondroyd, A. Bang Oturai, A. Pagnamenta, S. Y. Patel, M. Percy, N. Petousi, P. Piazza, S. E. Piret, G. Polanco-Echeverry, N. Popitsch, F. Powrie, C. Pugh, L. Quek, P. A. Robbins, K. Robson, A. Russo, N. Sahgal, P. A. van Schouwenburg, A. Schuh, E. Silverman, A. Simmons, P. S. Sørensen, E. Sweeney, J. Taylor, R. V. Thakker, I. Tomlinson, A. Trebes, S. R. F. Twigg, H. H. Uhlig, P. Vyas, T. Vyse, S. A. Wall, H. Watkins, M. P. Whyte, L. Witty, B. Wright, C. Yau, D. Buck, S. Humphray, P. J. Ratcliffe, J. I. Bell, A. O. M. Wilkie, D. Bentley, P. Donnelly, G. McVean, Factors influencing success of clinical genome sequencing across a broad spectrum of disorders. *Nat. Genet.* **47**, 717–726 (2015).
4. J. X. Chong, K. J. Buckingham, S. N. Jhangiani, C. Boehm, N. Sobreira, J. D. Smith, T. M. Harrell, M. J. McMillin, W. Wiszniewski, T. Gambin, Z. H. Coban Akdemir, K. Doheny, A. F. Scott, D. Avramopoulos, A. Chakravarti, J. Hoover-Fong, D. Mathews, P. D. Witmer, H. Ling, K. Hetrick, L. Watkins, K. E. Patterson, F. Reinier, E. Blue, D. Muzny, M. Kircher, K. Bilguvar, F. López-Giráldez, V. R. Sutton, H. K. Tabor, S. M. Leal, M. Gunel, S. Mane, R. A. Gibbs, E. Boerwinkle, A. Hamosh, J. Shendure, J. R. Lupski, R. P. Lifton, D. Valle, D. A. Nickerson; Centers for Mendelian Genomics, M. J. Bamshad, The genetic basis of Mendelian phenotypes: Discoveries, challenges, and opportunities. *Am. J. Hum. Genet.* **97**, 199–215 (2015).
5. D. MacArthur, T. A. Manolio, D. P. Dimmock, H. L. Rehm, J. Shendure, G. R. Abecasis, D. R. Adams, R. B. Altman, S. E. Antonarakis, E. A. Ashley, J. C. Barrett, L. G. Biesecker, D. F. Conrad, G. M. Cooper, N. J. Cox, M. J. Daly, M. B. Gerstein, D. B. Goldstein, J. N. Hirschhorn, S. M. Leal, L. A. Pennacchio, J. A. Stamatoyannopoulos, S. R. Sunyaev, D. Valle, B. F. Voight, W. Winkler, C. Gunter, Guidelines for investigating causality of sequence variants in human disease. *Nature* **508**, 469–476 (2014).
6. D. B. Goldstein, A. Allen, J. Keebler, E. H. Margulies, S. Petrou, S. Petrovski, S. Sunyaev, Sequencing studies in human genetics: Design and interpretation. *Nat. Rev. Genet.* **14**, 460–470 (2013).
7. M. Lek, D. MacArthur, The challenge of next generation sequencing in the context of neuromuscular diseases. *J. Neuromuscul. Dis.* **1**, 135–149 (2014).
8. Z. Wang, M. Gerstein, M. Snyder, RNA-seq: A revolutionary tool for transcriptomics. *Nat. Rev. Genet.* **10**, 57–63 (2009).
9. S. A. Byron, K. R. Van Keuren-Jensen, D. M. Engelthaler, J. D. Carpten, D. W. Craig, Translating RNA sequencing into clinical diagnostics: Opportunities and challenges. *Nat. Rev. Genet.* **17**, 257–271 (2016).
10. J. Tazi, N. Bakkour, S. Stamm, Alternative splicing and disease. *Biochim. Biophys. Acta* **1792**, 14–26 (2009).
11. P. Colapietro, P. Colapietro, C. Gervasini, F. Natacci, L. Rossi, P. Riva, L. Larizza, *NF1* exon 7 skipping and sequence alterations in exonic splice enhancers (ESEs) in a neurofibromatosis 1 patient. *Hum. Genet.* **113**, 551–554 (2003).
12. C. F. Morel, M. A. Thomas, H. Cao, C. H. O’Neil, J. G. Pickering, W. D. Foulkes, R. A. Hegele, A *LMNA* splicing mutation in two sisters with severe Dunnigan-type familial partial lipodystrophy type 2. *J. Clin. Endocrinol. Metabol.* **91**, 2689–2695 (2006).
13. M. Eriksson, W. Ted Brown, L. B. Gordon, M. W. Glynn, J. Singer, L. Scott, M. R. Erdos, C. M. Robbins, T. Y. Moses, P. Berglund, A. Dutra, E. Pak, S. Durkin, A. B. Csoka, M. Boehnke, T. W. Glover, F. S. Collins, Recurrent *de novo* point mutations in lamin A cause Hutchinson–Gilford progeria syndrome. *Nature* **423**, 293–298 (2003).
14. H. Gonorazky, M. Liang, B. Cummings, M. Lek, J. Micallef, C. Hawkins, R. Basran, R. Cohn, M. D. Wilson, D. MacArthur, C. R. Marshall, P. N. Ray, J. J. Dowling, RNAseq analysis for the diagnosis of muscular dystrophy. *Ann. Clin. Transl. Neurol.* **3**, 55–60 (2016).
15. K. Wang, C. Kim, J. Bradfield, Y. Guo, E. Toskala, F. G. Otieno, C. Hou, K. Thomas, C. Cardinale, G. J. Lyon, R. Golhar, H. Hakonarson, Whole-genome DNA/RNA sequencing identifies truncating mutations in *RBCK1* in a novel Mendelian disease with neuromuscular and cardiac involvement. *Genome Med.* **5**, 67 (2013).
16. H. Jung, D. Lee, J. Lee, D. Park, Y. Jeong Kim, W.-Y. Park, D. Hong, P. J. Park, E. Lee, Intron retention is a widespread mechanism of tumor-suppressor inactivation. *Nat. Genet.* **47**, 1242–1248 (2015).
17. Y. I. Li, B. van de Geijn, A. Raj, D. A. Knowles, A. A. Petti, D. Golan, Y. Gilad, J. K. Pritchard, RNA splicing is a primary link between genetic variation and disease. *Science* **352**, 600–604 (2016).
18. M. Melé, P. G. Ferreira, F. Reverter, D. S. DeLuca, J. Monlong, M. Sammeth, T. R. Young, J. M. Goldmann, D. D. Pervouchine, T. J. Sullivan, R. Johnson, A. V. Segrè, S. Djebali, A. Niarchou; GTEx Consortium, F. A. Wright, T. Lappalainen, M. Calvo, G. Getz, E. T. Dermitzakis, K. G. Ardlie, R. Guigó, The human transcriptome across tissues and individuals. *Science* **348**, 660–665 (2015).
19. E. T. Wang, R. Sandberg, S. Luo, I. Khrebtkova, L. Zhang, C. Mayr, S. F. Kingsmore, G. P. Schroth, C. B. Burge, Alternative isoform regulation in human tissue transcriptomes. *Nature* **456**, 470–476 (2008).
20. The GTEx Consortium, The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue gene regulation in humans. *Science* **348**, 648–660 (2015).
21. C. G. Bönnemann, C. H. Wang, S. Quijano-Roy, N. Deconinck, E. Bertini, A. Ferreira, F. Muntoni, C. Sewry, K. Bérout, K. D. Mathews, S. A. Moore, J. Bellini, A. Rutkowski, K. N. North; Members of the International Standard of Care Committee for Congenital Muscular Dystrophies, Diagnostic approach to the congenital muscular dystrophies. *Neuromuscul. Disord.* **24**, 289–311 (2014).
22. C. M. McDonald, Clinical approach to the diagnostic evaluation of hereditary and acquired neuromuscular diseases. *Phys. Med. Rehabil. Clin. N. Am.* **23**, 495–563 (2012).
23. S. Richards, N. Aziz, S. Bale, D. Bick, S. Das, J. Gastier-Foster, W. W. Grody, M. Hegde, E. Lyon, E. Spector, K. Voelkerding, H. L. Rehm; ACMG Laboratory Quality Assurance Committee, 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).
24. R. L. Begay, S. Graw, G. Sinagra, M. Merlo, D. Slavov, K. Gowan, K. L. Jones, G. Barbat, A. Spezzacatene, F. Brun, A. Di Lenarda, J. E. Smith, H. L. Granzier, L. Mestroni, M. Taylor; Familial Cardiomyopathy Registry, Role of titin missense variants in dilated cardiomyopathy. *J. Am. Heart Assoc.* **4**, e002645 (2015).
25. X. Roca, A. R. Krainer, I. C. Eperon, Pick one, but be quick: 5’ splice sites and the problems of too many choices. *Genes Dev.* **27**, 129–144 (2013).
26. M. A. Rivas, M. Pirinen, D. F. Conrad, M. Lek, E. K. Tsang, K. J. Karczewski, J. B. Maller, K. R. Kukurba, D. S. DeLuca, M. Fromer, P. G. Ferreira, K. S. Smith, R. Zhang, F. Zhao, E. Banks, R. Poplin, D. M. Ruderfer, S. M. Purcell, T. Tukiainen, E. V. Minikel, P. D. Stenson, D. N. Cooper, K. H. Huang, T. J. Sullivan, J. Nedzel; GTEx Consortium, Geuvadis Consortium, C. D. Bustamante, J. Billy Li, M. J. Daly, R. Guigó, P. Donnelly, K. Ardlie, M. Sammeth, E. T. Dermitzakis, M. I. McCarthy, S. B. Montgomery, T. Lappalainen, D. G. MacArthur, Effect of predicted protein-truncating genetic variants on the human transcriptome. *Science* **348**, 666–669 (2015).
27. K. Kiiski, V.-L. Lehtokari, A. Löytynoja, A. Ahlsten, J. Laitila, C. Wallgren-Pettersson, K. Pelin, A recursion copy number variation of the *NEB* triplicate region: Only revealed by the targeted nemaline myopathy CGH array. *Eur. J. Hum. Genet.* **24**, 574–580 (2015).
28. M.-L. Bang, T. Centner, F. Fornoff, A. J. Geach, M. Gotthardt, M. McNabb, C. C. Witt, D. Labeit, C. C. Gregorio, H. Granzier, S. Labeit, The complete gene sequence of titin, expression of an unusual ≈700-kDa titin isoform, and its interaction with obscurin identify a novel Z-line to I-band linking system. *Circ. Res.* **89**, 1065–1072 (2001).
29. The 1000 Genomes Project Consortium, A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
30. R. J. Butterfield, A. R. Foley, J. Dastgir, S. Asman, D. M. Dunn, Y. Zou, Y. Hu, S. Donkervoort, K. M. Flanagan, K. J. Swoboda, T. L. Winder, R. B. Weiss, C. G. Bönnemann, Position of glycine substitutions in the triple helix of *COL6A1*, *COL6A2*, and *COL6A3* is correlated with severity and mode of inheritance in collagen VI myopathies. *Hum. Mutat.* **34**, 1558–1567 (2013).
31. A. B. Spurdle, F. J. Couch, F. B. Hogervorst, P. Radice, O. M. Sinilnikova; IARC Unclassified Genetic Variants Working Group, Prediction and assessment of splicing alterations: Implications for clinical testing. *Hum. Mutat.* **29**, 1304–1313 (2008).
32. H. Duzkale, J. Shen, H. McLaughlin, A. Alfares, M. A. Kelly, T. J. Pugh, B. H. Funke, H. L. Rehm, M. S. Lebo, A systematic approach to assessing the clinical significance of genetic variants. *Clin. Genet.* **84**, 453–463 (2013).
33. L. S. Kremer, D. M. Bader, C. Mertes, R. Kopajtich, G. Pichler, A. Iuso, T. B. Haack, E. Graf, T. Schwarzmayr, C. Terrile, E. Konafikova, B. Repp, G. Kastenmüller, J. Adamski, P. Lichtner, C. Leonhardt, B. Funalot, A. Donati, V. Tiranti, A. Lombes, C. Jardel, D. Gläser, R. W. Taylor, D. Ghezzi, J. A. Mayr, A. Rötig, P. Freisinger, F. Distelmaier, T. M. Strom, T. Meitinger, J. Gagneur, Genetic diagnosis of Mendelian disorders via RNA sequencing. *bioRxiv* **2016**, 10.1101/066738 (2016).
34. D. S. DeLuca, J. Z. Levin, A. Sivachenko, T. Fennell, M.-D. Nazaire, C. Williams, M. Reich, W. Winkler, G. Getz, RNA-SeQC: RNA-seq metrics for quality control and process optimization. *Bioinformatics* **28**, 1530–1532 (2012).
35. S. Schafer, K. Miao, C. C. Benson, M. Heinig, S. A. Cook, N. Hubner, Alternative splicing signatures in RNA-seq data: Percent spliced in (PSI). *Curr. Protoc. Hum. Genet.* **87**, 11.16. 1–11.16. 14 (2015).

36. S. M. Purcell, J. L. Moran, M. Fromer, D. Ruderfer, N. Solovieff, P. Roussos, C. O'Dushlaine, K. Chambert, S. E. Bergen, A. Kähler, L. Duncan, E. Stahl, G. Genovese, E. Fernández, M. O. Collins, N. H. Komiya, J. S. Choudhary, P. K. E. Magnusson, E. Banks, K. Shakir, K. Garimella, T. Fennell, M. DePristo, S. G. N. Grant, S. J. Haggarty, S. Gabriel, E. M. Scolnick, E. S. Lander, C. M. Hultman, P. F. Sullivan, S. A. McCarroll, P. Sklar, A polygenic burden of rare disruptive mutations in schizophrenia. *Nature* **506**, 185–190 (2014).
37. M. Lek, K. J. Karczewski, E. V. Minikel, K. E. Samocha, E. Banks, T. Fennell, A. H. O'Donnell-Luria, J. S. Ware, A. J. Hill, B. B. Cummings, T. Tukiainen, D. P. Birnbaum, J. A. Kosmicki, L. E. Duncan, K. Estrada, F. Zhao, J. Zou, E. Pierce-Hoffman, J. Berghout, D. N. Cooper, N. Deflaux, M. DePristo, R. Do, J. Flannick, M. Fromer, L. Gauthier, J. Goldstein, N. Gupta, D. Howrigan, A. Kiezun, M. I. Kurki, A. Levy Moonshine, P. Natarajan, L. Orozco, G. M. Peloso, R. Poplin, M. A. Rivas, V. Ruano-Rubio, S. A. Rose, D. M. Ruderfer, K. Shakir, P. D. Stenson, C. Stevens, B. P. Thomas, G. Tiao, M. T. Tusie-Luna, B. Weisburd, H.-H. Won, D. Yu, D. M. Altshuler, D. Ardisino, M. Boehnke, J. Danesh, S. Donnelly, R. Elosua, J. C. Florez, S. B. Gabriel, G. Getz, S. J. Glatt, C. M. Hultman, S. Kathiresan, M. Laakso, S. McCarroll, M. I. McCarthy, D. McGovern, R. McPherson, B. M. Neale, A. Palotie, S. M. Purcell, D. Saleheen, J. M. Scharf, P. Sklar, P. F. Sullivan, J. Tuomilehto, M. T. Tsuang, H. C. Watkins, J. G. Wilson, M. J. Daly, D. G. MacArthur; Exome Aggregation Consortium, Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285–291 (2016).
38. S. E. Castel, A. Levy-Moonshine, P. Mohammadi, E. Banks, T. Lappalainen, Tools and best practices for data processing in allelic expression analysis. *Genome Biol.* **16**, 1–12 (2015).
39. G. Yeo, C. B. Burge, Maximum entropy modeling of short sequence motifs with applications to RNA splicing signals. *J. Comput. Biol.* **11**, 377–394 (2004).
40. M. G. Reese, F. H. Eckman, D. Kulp, D. Haussler, Improved splice site detection in Genie. *J. Comput. Biol.* **4**, 311–323 (1997).
41. M. Pertea, X. Lin, S. L. Salzberg, GeneSplicer: A new computational method for splice site prediction. *Nucleic Acids Res.* **29**, 1185–1190 (2001).
42. F.-O. Desmet, D. Hamroun, M. Lalande, G. Colod-Bérout, M. Claustres, C. Bérout, Human Splicing Finder: An online bioinformatics tool to predict splicing signals. *Nucleic Acids Res.* **37**, e67 (2009).
43. B. Mersch, A. Geppert, S. Suhai, A. Hotz-Wagenblatt, Automatic detection of exonic splicing enhancers (ESEs) using SVMs. *BMC Bioinformatics* **9**, 369 (2008).
44. L. A. Chasin, Searching for splicing motifs. *Adv. Exp. Med. Biol.* **623**, 85–106 (2008).
45. S. T. Cooper, H. P. Lo, K. N. North, Single section Western blot Improving the molecular diagnosis of the muscular dystrophies. *Neurology* **61**, 93–97 (2003).
46. Y.-C. Chan, H.-Q. Tong, A. H. Beggs, L. M. Kunkel, Human skeletal muscle-specific α -actinin-2 and -3 isoforms form homodimers and heterodimers in vitro and in vivo. *Biochem. Biophys. Res. Commun.* **248**, 134–139 (1998).

Acknowledgments: Sequencing and analysis were provided by the Broad Institute of Massachusetts Institute of Technology (MIT) and Harvard Center for Mendelian Genomics (Broad CMG). We thank H. Brooks, D. Sookiasian, M. E. Leach, D. Ezzi, J. Dastgir, A. Rutkowski, C. Grosman, C. Konermans, S. Ceulemans, M.-L. Chu, E. Moran, and K. Matthews for sample collection and quality control. We also thank C. Miceli, S. Nelson, V. Rusu, and D. Altshuler for sharing control cell lines and plasmids. **Funding:** This project was supported by funding from the Broad Institute's Broadlight and Broadnext10 programs. B.B.C. is supported by the NIH GM096911 training grant. T.T. is supported by the Academy of Finland, the Finnish Cultural Foundation, the Orion-Farmos Research Foundation, and the Emil Aaltonen Foundation. M.L. is supported by the Australian NHMRC (National Health and Medical Research Council) CJ Martin Fellowship, the Australian American Association Sir Keith Murdoch Fellowship, and a Muscular Dystrophy Association/American Association of Neuromuscular and Electrodiagnostic Medicine (MDA/AANEM) development grant. L.B.W., S.A.S., N.G.L., N.F.C., K.N.N., and E.C.O. are supported by the NHMRC of Australia (1080587, 1075451, 1002147, 1113531, 1022707, 1031893, and 1090428). K.J.K. is supported by a National Institute of General Medical Sciences (NIGMS) fellowship grant (F32GM115208). A.H.O.-L. is supported by an NIGMS fellowship grant (4T32GM007748). A.H.B. is supported by the NIH R01 HD075802 and R01 AR044345 and by MDA383249 from the Muscular Dystrophy Association. P.B.K., E.E., and H.K.M. are supported by NIH R01NS080929. J.J.D. is supported in part by funding from Genome Canada (a Disruptive Innovations in Genomics grant). Funding relevant to this research includes fellowship support of S.T.C. and a project grant supporting an Australian-wide program about gene discovery in inherited neuromuscular disorders performed in collaboration with D.G.M. [NHMRC APP1048816 (2013–2017) and NHMRC APP1080587 (2015–2019)]. The Broad CMG was funded by the National Human Genome Research Institute (NHGRI), the National Eye Institute, and the National Heart, Lung, and Blood Institute (NHLBI) grant UM1 HG008900 to D.G.M. and H. Rehm. The GTEx project was supported by the Common Fund of the Office of the Director of the NIH (<http://commonfund.nih.gov/GTEx>). Additional funds were provided by the National Cancer Institute (NCI), NHGRI, NHLBI, National Institute on Drug Abuse (NIDA), National Institute of Mental Health (NIMH), and National Institute of Neurological Disorders and Stroke (NINDS). Donors were enrolled at Biospecimen Source Sites that were funded by NCI/Science Applications

International Corporation (SAIC)–Frederick Inc. (SAIC-F) subcontracts to the National Disease Research Interchange (10XS170) and the Roswell Park Cancer Institute (10XS171). The Laboratory, Data Analysis, and Coordinating Center (LDACC) was funded through a contract (HHSN268201000029C) to the Broad Institute Inc. Biorepository operations were funded through an SAIC-F subcontract to the Van Andel Institute (10ST1035). Additional data repository and project management were provided by SAIC-F (HHSN261200800001E). The Brain Bank was supported by a supplement to the University of Miami grant DA006227.

Author contributions: B.B.C., T.T., and D.G.M. conceived and designed the experiments. B.B.C. and T.T. analyzed the RNA-seq data. J.L.M., Y.H., A.B., and M.R.D. designed and performed validation experiments. B.B.C., M.L., S.D., A.R.F., L.B.W., S.A.S., G.L.O., H.M.R., E.C.O., R.G., S.T.C., and C.G.B. analyzed the exome and whole-genome data. S.D., A.R.F., V.B., L.B.W., S.A.S., G.L.O., E.E., H.M.R., A.S., H.G., K.C., E.C.O., R.G., N.G.L., A.T., A.H.B., P.B.K., K.N.N., V.S., J.J.D., F.M., N.F.C., S.T.C., and C.G.B. provided patient samples and clinical information. The Broad CMG and GTEx provided sequencing support for patient and control DNA sequencing and RNA-seq. F.Z., B.W., K.J.K., A.H.O.-L., D.B., and H.J. contributed reagents, materials, and analysis tools. J.L.M., T.T., M.L., S.D., A.R.F., V.B., L.B.W., S.A.S., K.J.K., A.H.O.-L., E.C.O., N.G.L., A.T., J.J.D., C.G.B., and S.T.C. critically evaluated the manuscript. B.B.C. and D.G.M. wrote the manuscript. **Competing interests:** C.G.B., V.B., D.G.M., M.L., B.B.C., and S. Wilton are inventors on U.S. Provisional Patent Application no. 62/358,482, which covers "Diagnosing COL6-related disorders and methods for treating same," and was filed on 5 July 2016 by NINDS. D.G.M. is a founder and owns stock in Goldfinch Biopharma, but this work is unrelated to the company. All other authors declare that they have no competing interests. **Data and materials availability:** Patient sequencing data generated as part of this study were deposited in dbGaP under accession ID phs000655.v3.p1. GTEx transcriptome sequencing data can be obtained from dbGaP under accession ID phs000424.v6.p1. Code for splice junction discovery, normalization, and filtering is available on <https://github.com/berylc/MendelianRNA-seq>. List of OMIM and neuromuscular disease genes used for splice detection and allele-specific expression analysis can be found at <https://github.com/macarthur-lab/omim> and <https://github.com/berylc/MendelianRNA-seq>, respectively.

Members of the GTEx Consortium

LDACC–Analysis Working Group (AWG): Kristin G. Ardlie,¹ Gad Getz,^{1,2} Ellen T. Gelfand,¹ Ayellet V. Segrè,¹ François Aguet,¹ Timothy J. Sullivan,¹ Xiao Li,¹ Jared L. Nedzel,¹ Casandra A. Trowbridge,¹ Daniel G. MacArthur,^{1,3} Monkol Lek,^{1,3} Taru Tukiainen,^{3,4} Kane Hadley,⁴ Katherine H. Huang,⁴ Michael S. Noble,⁴ Duyen T. Nguyen,⁴ Beryl B. Cummings,^{3,4} **Funded Statistical Methods groups–AWG:** Andrew B. Nobel,⁵ Fred A. Wright,⁶ Andrey A. Shabalov,⁷ John J. Palowitch,⁸ Yi-Hui Zhou,⁹ Emmanouil T. Dermizakis,^{10,11,12} Mark I. McCarthy,^{13,14,15} Anthony J. Payne,¹³ Tuuli Lappalainen,^{16,17} Stephane Castel,^{16,17} Sarah Kim-Hellmuth,^{16,17} Pejman Mohammadi,^{16,17} Alexis Battle,¹⁸ Princy Parsana,¹⁸ Sara Mostafavi,¹⁹ Andrew Brown,^{10,11,12} Halit Ongen,^{10,11,12} Olivier Delaneau,^{10,11,12} Nikolaos Panousis,^{10,11,12} Cedric Howald,^{10,11,12} Martijn van de Bunt,^{13,14} Roderic Guigo,^{20,21,22} Jean Monlong,^{20,21,23} Ferran Reverter,^{20,24} Diego Garrido,^{20,21} Manuel Munoz,^{20,21} Gireesh Bogu,^{20,21} Reza Sodaei,^{20,21} Panagiotis Papasaiaks,^{20,21} Anne W. Ndungu,¹³ Stephen B. Montgomery,²⁵ Xin Li,²⁵ Laure Fresard,²⁵ Joe R. Davis,²⁵ Emily K. Tsang,^{25,26} Zachary Zappala,²⁵ Nathan S. Abell,²⁵ Michael J. Gloudemans,^{25,26} Boxiang Liu,^{25,27} Farhan N. Damani,²⁸ Ashis Saha,²⁸ Yungil Kim,¹⁸ Benjamin J. Strobe,²⁹ Yuan He,²⁹ Matthew Stephens,^{30,31} Jonathan K. Pritchard,^{30,32,33} Xiaoquan Wen,³⁴ Sarah Urbat,³⁰ Nancy J. Cox,^{35,36} Dan L. Nicolae,³⁷ Eric R. Gamazon,^{35,36} Hae Kyung Im,³⁸ Christopher D. Brown,³⁹ Barbara E. Engelhardt,⁴⁰ YoSon Park,³⁹ Brian Jo,⁴¹ Ian C. McDowell,⁴² Ariel Gewirtz,⁴¹ Genna Gliner,⁴³ Don Conrad,^{44,45} Ira Hall,^{46,47,48} Colby Chiang,⁴⁶ Alexandra Scott,⁴⁶ Chiara Sabatti,⁴⁹ Eleazar Eskin,⁵⁰ Christine Peterson,⁵¹ Farhad Hormozdizari,⁵² Eun Yong Kang,⁵² Serghei Mangul,⁵² Buhm Han,⁵³ Jae Hoon Sul,⁵⁴ **Enhancing GTEx funded group:** Andrew P. Feinberg,⁵⁵ Lindsay F. Rizzardi,⁵⁶ Kasper D. Hansen,⁵⁷ Peter Hickey,⁵⁹ Joshua Akey,⁵⁹ Manolis Kellis,⁶⁰ Jin Billy Li,⁶¹ Michael Snyder,⁶¹ Hua Tang,⁶¹ Lihua Jiang,⁶¹ Shin Lin,^{61,62} Barbara E. Stranger,⁶³ Marian Fernando,⁶⁴ Meritxell Oliva,⁶⁴ John Stamatoyannopoulos,⁶⁵ Rajinder Kaul,⁶⁵ Jessica Halow,⁶⁵ Richard Sandstrom,⁶⁵ Eric Haugen,⁶⁵ Audra Johnson,⁶⁵ Kristen Lee,⁶⁵ Daniel Bates,⁶⁵ Morgan Diegel,⁶⁵ Brandon L. Pierce,⁶⁶ Lin Chen,⁶⁶ Muhammad G. Kibriya,⁶⁶ Farzana Jasmine,⁶⁶ Jennifer Doherty,⁶⁷ Kathryn Demanelis,⁶⁶ Stephen B. Montgomery,²⁵ Emily K. Tsang,²⁵ Kevin S. Smith,²⁵ Qin Li,⁶¹ Rui Zhang,⁶¹ **National Institutes of Health (NIH) Common Fund:** Concepcion R. Nierras,⁶⁸ **NIH/NCI:** Helen M. Moore,⁶⁹ Abhi Rao,⁶⁹ Ping Guan,⁶⁹ Jimmie B. Vaughn,⁶⁹ Philip A. Branton,⁶⁹ Latarsha J. Carithers,⁷⁰ **NIH/NHGRI:** Simona Volpi,⁷¹ Jeffery P. Struwing,⁷¹ Casey G. Martin,⁷¹ Lockhart C. Nicole,⁷² **NIH/NIMH:** Susan E. Koester,⁷² Anjene M. Addington,⁷² **NIH/NIDA:** A. Roger. Little,⁷³ **Biospecimen Collection Source Site–National Disease Research Interchange:** William F. Leinweber,⁷⁴ Jeffrey A. Thomas,⁷⁴ Gene Kopen,⁷⁴ Alisa McDonald,⁷⁴ Bernadette Mestichelli,⁷⁴ Saboor Shad,⁷⁴ John T. Lonsdale,⁷⁴ Michael Salvatore,⁷⁴ Richard Hasz,⁷⁵ Gary Walters,⁷⁶ Mark Johnson,⁷⁶ Michael Washington,⁷⁶ Lori E. Brigham,⁷⁷ Christopher Johns,⁷⁸ Joseph Wheeler,⁷⁸ Brian Roe,⁷⁹ Marcus Hunter,⁷⁹ Kevin Myer,⁷⁹ **Biospecimen Collection Source Site–Roswell Park Cancer Institute:** Barbara A. Foster,⁸⁰ Michael T. Moser,⁸⁰ Ellen Karasik,⁸⁰ Bryan M. Gillard,⁸⁰ Rachna Kumar,⁸⁰ Jason Bridge,⁸¹ Mark Miklos,⁸¹ **Biospecimen Core Resource–Van Andel Research Institute:** Scott D. Jewell,⁸² Daniel C. Rohrer,⁸² Dana Valley,⁸² Robert G. Montroy,⁸² **Brain Bank Repository–University**

of Miami: Deborah C. Mash,⁸³ David A. Davis,⁸⁴ **Leidos Biomedical Project Management:** Anita H. Undale,⁸⁵ Anna M. Smith,⁸⁶ David E. Tabor,⁸⁶ Nancy V. Roche,⁸⁶ Jeffrey A. McLean,⁸⁶ Negin Vatanian,⁸⁶ Karna L. Robinson,⁸⁶ Leslie Sobin,⁸⁶ Mary E. Barcus,⁸⁷ Kimberly M. Valentino,⁸⁶ Liquan Qi,⁸⁶ Stephen Hunter,⁸⁶ Pushpa Hariharan,⁸⁶ Shilpi Singh,⁸⁶ Ki Sung Um,⁸⁶ Takunda Matose,⁸⁶ Maria M. Tomadzewski,⁸⁶ **Ethical, Legal, and Social Implications Study:** Laura A. Siminoff,⁸⁸ Heather M. Traino,⁸⁹ Maghboeba Mosavel,⁹⁰ Laura K. Barker,⁹¹ **Genome Browser Data Integration, and Visualization–European Bioinformatics Institute:** Daniel R. Zerbino,⁹² Thomas Juettmann,⁹² Kieron Taylor,⁹² Magali Ruffier,⁹² Dan Sheppard,⁹² Steven Trevanion,⁹² Paul Flicek,⁹² **Genome Browser Data Integration and Visualization–Genomics Institute, University of California, Santa Cruz:** W. James Kent,⁹³ Kate R. Rosenbloom,⁹³ Maximilian Haessler,⁹³ Christopher M. Lee,⁹³ Benedict Paten,⁹³ John Vivan,⁹³ Jingchun Zhu,⁹³ Mary Goldman,⁹³ Brian Craft,⁹³ **Other members of the AWG:** Gen Li,⁹⁴ Pedro G. Ferreira,^{95,96} Esti Yeger-Lotem,^{97,98} Matthew T. Maurano,⁹⁹ Ruth Barshir,⁹⁷ Omer Basha,⁹⁷ Hualin S. Xi,¹⁰⁰ Jie Quan,¹⁰⁰ Michael Sammeth,¹⁰¹ Judith B. Zaugg¹⁰²

¹Broad Institute of MIT and Harvard University, Cambridge, MA 02142, USA. ²Massachusetts General Hospital Cancer Center and Department of Pathology, Massachusetts General Hospital and Harvard Medical School, Boston, MA 02114, USA. ³Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA 02114, USA. ⁴Broad Institute of MIT and Harvard University, Cambridge, MA 02142, USA. ⁵Department of Statistics and Operations Research and Department of Biostatistics, University of North Carolina at Chapel Hill, Chapel Hill, NC 27599–3260, USA. ⁶Bioinformatics Research Center and Departments of Statistics and Biological Sciences, North Carolina State University, Raleigh, NC 27695, USA. ⁷Center for Biomarker Research and Personalized Medicine, Virginia Commonwealth University, Richmond, VA 23298–0581, USA. ⁸Department of Statistics and Operations Research, University of North Carolina at Chapel Hill, Chapel Hill, NC 27599–3260, USA. ⁹Bioinformatics Research Center and Department of Biological Sciences, North Carolina State University, Raleigh, NC 27695, USA. ¹⁰Department of Genetic Medicine and Development, University of Geneva Medical School, 1211 Geneva, Switzerland. ¹¹Institute for Genetics and Genomics in Geneva (iGE3), University of Geneva, 1211 Geneva, Switzerland. ¹²Swiss Institute of Bioinformatics, 1211 Geneva, Switzerland. ¹³Wellcome Trust Centre for Human Genetics Research, Nuffield Department of Clinical Medicine, University of Oxford, Oxford OX3 7BN, U.K. ¹⁴Oxford Centre for Diabetes, Endocrinology and Metabolism, Churchill Hospital, University of Oxford, Oxford OX3 7LE, U.K. ¹⁵Oxford National Institute for Health Research Biomedical Research Centre, Churchill Hospital, Oxford OX3 7LJ, U.K. ¹⁶New York Genome Center, 101 Avenue of the Americas, New York, NY 10013, USA. ¹⁷Department of Systems Biology, Columbia University Medical Center, New York, NY 10032, USA. ¹⁸Department of Computer Science, Johns Hopkins University, Baltimore, MD 21218, USA. ¹⁹Department of Computer Science, Stanford University, Stanford, CA 94305, USA. ²⁰Center for Genomic Regulation, Barcelona, Catalonia, Spain. ²¹Universitat Pompeu Fabra, 08003 Barcelona, Catalonia, Spain. ²²Institut Hospital del Mar d'Investigacions Mèdiques, 08003 Barcelona, Spain. ²³Department of Human Genetics, McGill University, Montréal, Québec, Canada. ²⁴Universitat de Barcelona, 08028 Barcelona, Catalonia, Spain. ²⁵Departments of Genetics and Pathology, Stanford University, Stanford, CA 94305, USA. ²⁶Biomedical Informatics Program, Stanford University, Stanford, CA 94305, USA. ²⁷Department of Biology, Stanford University, Stanford, CA 94305, USA. ²⁸Department of Computer Science, Johns Hopkins University, Baltimore, MD 21218, USA. ²⁹Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD 21218, USA. ³⁰Department of Human Genetics, University of Chicago, Chicago, IL 60637, USA. ³¹Department of Statistics, University of Chicago, 5734 South University Avenue, Chicago, IL 60637, USA. ³²Departments of Genetics and Biology, Stanford University, Stanford, CA 94305, USA. ³³Howard Hughes Medical Institute, Chevy Chase, MD 10032, USA. ³⁴Department of Biostatistics, University of Michigan, 1415 Washington Heights, Ann Arbor, MI 48109, USA. ³⁵Division of Genetic Medicine, Department of Medicine, Vanderbilt University Medical Center, Nashville, TN 37232, USA. ³⁶Department of Clinical Epidemiology, Biostatistics and Bioinformatics and Department of Psychiatry, Academic Medical Center, University of Amsterdam, Meibergdreef 9, 1105 AZ Amsterdam, Netherlands. ³⁷Section of Genetic Medicine, Department of Medicine, Department of Statistics, and Department of Human Genetics, University of Chicago, 900 East 57th Street KCBD 3220, Chicago, IL 60637, USA. ³⁸Section of Genetic Medicine, Department of Medicine, University of Chicago, 900 East 57th Street KCBD 3220, Chicago, IL 60637, USA. ³⁹Department of Genetics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA 19104, USA. ⁴⁰Department of Computer Science, Center for Statistics and Machine Learning, Princeton University, 35 Olden Street, Princeton, NJ 08540, USA. ⁴¹Lewis-Sigler Institute, Princeton University, Princeton, NJ 08540, USA. ⁴²Computational Biology and Bioinformatics Graduate Program, Duke University, Durham, NC 27708, USA. ⁴³Department of Operations Research and Financial Engineering, Princeton University, Princeton, NJ 08540, USA. ⁴⁴Department of Genetics, Washington University School of Medicine, St. Louis, MO 63108, USA. ⁴⁵Department of Pathology and Immunology, Washington University School of Medicine, St. Louis, MO 63108, USA. ⁴⁶McDonnell Genome Institute, Washington University School of Medicine, Saint Louis, MO 63108, USA. ⁴⁷Department of Medicine, Washington University School of Medicine, Saint Louis, MO 63108, USA. ⁴⁸Department of Genetics, Washington University School of Medicine, Saint Louis, MO 63108, USA. ⁴⁹Departments of Biomedical Data Science and Statistics, Stanford University, Health Research and Policy Redwood building, Stanford, CA 94305–5404, USA. ⁵⁰Departments of Computer Science and

Human Genetics, University of California, Los Angeles, Los Angeles, CA 90095, USA. ⁵¹Department of Biostatistics, University of Texas MD Anderson Cancer Center, 1400 Pressler Street, Houston, TX 77030, USA. ⁵²Department of Computer Science, University of California, Los Angeles, Los Angeles, CA 90095, USA. ⁵³Department of Convergence Medicine, University of Ulsan College of Medicine, Asan Medical Center, Mugeo-dong, Nam-gu, Ulsan, Korea. ⁵⁴Department of Psychiatry and Biobehavioral Sciences, University of California, Los Angeles, Los Angeles, CA 90095, USA. ⁵⁵Center for Epigenetics, Johns Hopkins University School of Medicine, and Departments of Medicine, Biomedical Engineering, and Mental Health, Johns Hopkins University Schools of Medicine, Engineering, and Public Health, Baltimore, MD 21205, USA. ⁵⁶Center for Epigenetics, Johns Hopkins University School of Medicine, Baltimore, MD 21205, USA. ⁵⁷McKusick-Nathans Institute of Genetic Medicine, Center for Epigenetics, Johns Hopkins School of Medicine, and Department of Biostatistics, Bloomberg School of Public Health, Johns Hopkins University, Baltimore, MD 21205, USA. ⁵⁸Department of Biostatistics, Bloomberg School of Public Health, Johns Hopkins University, Baltimore, MD 21205, USA. ⁵⁹Department of Genome Sciences, University of Washington, Seattle, WA 98195, USA. ⁶⁰Computer Science and Artificial Intelligence Laboratory, MIT, Cambridge, MA 02139, USA. ⁶¹Department of Genetics, Stanford University, Stanford, CA 94305, USA. ⁶²Division of Cardiology, University of Washington, Seattle, WA 98195, USA. ⁶³Section of Genetic Medicine, Department of Medicine, Institute for Genomics and Systems Biology, Center for Data Intensive Science, University of Chicago, Chicago, IL 60637, USA. ⁶⁴Section of Genetic Medicine, Department of Medicine, Institute for Genomics and Systems Biology, University of Chicago, Chicago, IL 60637, USA. ⁶⁵Altius Institute for Biomedical Sciences, Seattle, WA 98121, USA. ⁶⁶Department of Public Health Sciences, University of Chicago, Chicago, IL 60637, USA. ⁶⁷Department of Epidemiology, Geisel School of Medicine at Dartmouth, Lebanon, NH 03756, USA. ⁶⁸Office of Strategic Coordination, Division of Program Coordination, Planning, and Strategic Initiatives, Rockville, MD 20852–9305, USA. ⁶⁹Biorepositories and Biospecimen Research Branch, Division of Cancer Treatment and Diagnosis, NCI, Bethesda, MD 20892, USA. ⁷⁰National Institute of Dental and Craniofacial Research, 6701 Democracy Boulevard, Bethesda, MD 20892, USA. ⁷¹Division of Genomic Medicine, NHGRI, Rockville, MD 20892, USA. ⁷²Division of Neuroscience and Basic Behavioral Science, NIMH, NIH, Bethesda, MD 20892, USA. ⁷³NIDA, NIH, U.S. Department of Health and Human Services, Bethesda, MD 20892, USA. ⁷⁴National Disease Research Interchange, Philadelphia, PA 19103, USA. ⁷⁵Gift of Life Donor Program, Philadelphia, PA 19103, USA. ⁷⁶LifeNet Health, Virginia Beach, VA 23453, USA. ⁷⁷Washington Regional Transplant Community, Annandale, VA 22003, USA. ⁷⁸Center for Organ Recovery and Education, Pittsburgh, PA 15238, USA. ⁷⁹LifeGift, Houston, TX 77054, USA. ⁸⁰Roswell Park Cancer Institute Pharmacology and Therapeutics, Buffalo, NY 14263, USA. ⁸¹Unyts, 110 Broadway, Buffalo, NY 14203, USA. ⁸²Van Andel Research Institute, Grand Rapids, MI 49503, USA. ⁸³Department of Neurology, Miller School of Medicine, University of Miami, Miami, FL 33136, USA. ⁸⁴Brain Endowment Bank, Miller School of Medicine, University of Miami, Miami, FL 33136, USA. ⁸⁵National Institute of Allergy and Infectious Diseases, NIH, 5601 Fishers Lane, Rockville, MD 20852, USA. ⁸⁶Biospecimen Research Group, Clinical Research Directorate, Leidos Biomedical Research Inc., Rockville, MD 20852, USA. ⁸⁷Frederick National Laboratory for Cancer Research, 8560 Progress Drive, Room C3021, Frederick, MD 21701, USA. ⁸⁸Temple University, Philadelphia, PA 19122, USA. ⁸⁹Temple University, Ritter Annex 9th Floor, 1301 Cecil B. Moore Avenue, Philadelphia, PA 19122, USA. ⁹⁰Virginia Commonwealth University, Richmond, VA 23219, USA. ⁹¹Temple University, Philadelphia, PA 19122, USA. ⁹²European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, U.K. ⁹³Genomics Institute, University of California, Santa Cruz, Santa Cruz, CA 95064, USA. ⁹⁴Department of Biostatistics, Mailman School of Public Health, Columbia University, New York, NY 10032, USA. ⁹⁵IS Instituto de Investigação e Inovação em Saúde, Universidade do Porto, Rua Alfredo Allen, 208, 4200–135 Porto, Portugal. ⁹⁶IPATIMUP–Institute of Molecular Pathology and Immunology, University of Porto, Rua Dr. Roberto Frias s/número, 4200–625 Porto, Portugal. ⁹⁷Ben-Gurion University of the Negev, Beer-Sheva, 84105 Israel. ⁹⁸National Institute for Biotechnology in the Negev, Beer-Sheva 84105, Israel. ⁹⁹Institute for Systems Genetics, New York University Langone Medical Center, New York, NY 10016, USA. ¹⁰⁰Computational Sciences, Pfizer Inc., 610 Main Street, Cambridge, MA 02140, USA. ¹⁰¹Institute of Biophysics Carlos Chagas Filho, Federal University of Rio de Janeiro (UFRJ), 21941902 Rio de Janeiro, Brazil. ¹⁰²European Molecular Biology Laboratory, Meyerhofstrasse 1, 69117 Heidelberg, Germany.

Submitted 9 January 2017

Accepted 29 March 2017

Published 19 April 2017

10.1126/scitranslmed.aal5209

Citation: B. B. Cummings, J. L. Marshall, T. Tukiainen, M. Lek, S. Donkervoort, A. R. Foley, V. Bolduc, L. B. Waddell, S. A. Sandaradura, G. L. O'Grady, E. Estrella, H. M. Reddy, F. Zhao, B. Weisburd, K. J. Karczewski, A. H. O'Donnell-Luria, D. Birnbaum, A. Sarkozy, Y. Hu, H. Gonorazky, K. Claeys, H. Joshi, A. Bournazos, E. C. Oates, R. Ghaoui, M. R. Davis, N. G. Laing, A. Topf, Genotype-Tissue Expression Consortium, P. B. Kang, A. H. Beggs, K. N. North, V. Straub, J. J. Dowling, F. Muntoni, N. F. Clarke, S. T. Cooper, C. G. Bönnemann, D. G. MacArthur, Improving genetic diagnosis in Mendelian disease with transcriptome sequencing. *Sci. Transl. Med.* 9, eaal5209 (2017).