

GREGoR: Accelerating Genomics for Rare Diseases

Moez Dawood^{1,2,3}, Ben Heavner⁴, Marsha M. Wheeler⁴, Rachel A. Ungar^{5,6,7}, Jonathan LoTempio⁸, Laurens Wiel^{5,6,9}, Seth Berger^{10,11,12}, Jonathan A. Bernstein¹³, Jessica X. Chong^{14,15}, Emmanuèle C. Délot⁸, Evan E. Eichler^{15,16,17}, Richard A. Gibbs^{1,2}, James R. Lupski^{1,2,18}, Ali Shojaie⁴, Michael E. Talkowski^{19,20,21,22,23}, Alex H. Wagner^{24,25,26}, Chia-Lin Wei¹⁶, Christopher Wellington²⁷, Matthew T. Wheeler⁹, GREGoR Partner Members, Claudia M. B. Carvalho²⁸, Casey A. Gifford^{5,13,29,30}, Susanne May⁴, Danny E. Miller^{15,16,31,32}, Heidi L. Rehm^{19,20}, Fritz J. Sedlazeck^{1,33}, Eric Vilain⁸, Anne O'Donnell-Luria^{19,20,34}, Jennifer E. Posey², Lisa H. Chadwick³⁵, Michael J. Bamshad^{14,15,36}, Stephen B. Montgomery^{5,6,37}, Genomics Research to Elucidate the Genetics of Rare Diseases (GREGoR) Consortium

¹Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX, USA.

²Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, TX, USA.

³Medical Scientist Training Program, Baylor College of Medicine, Houston, TX, USA.

⁴Department of Biostatistics, University of Washington, Seattle, WA, USA.

⁵Department of Genetics, School of Medicine, Stanford University, Stanford, CA, USA.

⁶Department of Pathology, School of Medicine, Stanford University, Stanford, CA, USA.

⁷Stanford Center for Biomedical Ethics, School of Medicine, Stanford University, Stanford, CA, USA.

⁸Institute for Clinical and Translational Science, University of California, Irvine, CA, USA.

⁹Division of Cardiovascular Medicine, School of Medicine, Stanford University, Stanford, CA, USA.

¹⁰Division of Genetics and Metabolism, Children's National Rare Disease Institute, Washington, DC, USA.

¹¹Center for Genetic Medicine Research, Children's National Rare Disease Institute, Washington, DC, USA.

¹²Department of Genomics and Precision Medicine, George Washington University, Washington, DC, USA.

¹³Department of Pediatrics, School of Medicine, Stanford University, Stanford, CA, USA.

¹⁴Department of Pediatrics, Division of Genetic Medicine, University of Washington, Seattle, WA, USA.

¹⁵Brotman Baty Institute for Precision Medicine, University of Washington, Seattle, WA, USA.

¹⁶Department of Genome Sciences, University of Washington, Seattle, WA, USA.

¹⁷Howard Hughes Medical Institute, University of Washington, Seattle, WA, USA.

¹⁸Department of Pediatrics, Baylor College of Medicine, Houston, TX, USA.

¹⁹Center for Genomic Medicine, Massachusetts General Hospital, Boston, MA, USA.

²⁰Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Boston, MA, USA.

²¹Department of Neurology, Massachusetts General Hospital and Harvard Medical School, Boston, MA, USA.

²²Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, MA, USA.

²³Program in Bioinformatics and Integrative Genomics, Harvard Medical School, Boston, MA, USA.

²⁴Steve and Cindy Rasmussen Institute for Genomic Medicine, Nationwide Children's Hospital, Columbus, OH, USA.

²⁵Department of Pediatrics, The Ohio State University College of Medicine, Columbus, OH, USA.

²⁶Department of Biomedical Informatics, The Ohio State University College of Medicine, Columbus, OH, USA.

²⁷Office of Genomic Data Science, National Human Genome Research Institute, Bethesda, MD, USA.

²⁸Pacific Northwest Research Institute, Seattle, WA, USA.

²⁹Basic Science and Engineering Initiative, Stanford Children's Health, Betty Irene Moore Children's Heart Center, Stanford, CA, USA.

³⁰Institute for Stem Cell Biology and Regenerative Medicine, School of Medicine, Stanford University, Stanford, CA, USA.

³¹Division of Genetic Medicine, Department of Pediatrics, University of Washington, Seattle, WA, USA.

³²Department of Laboratory Medicine and Pathology, University of Washington, Seattle, WA, USA.

³³Department of Computer Science, Rice University, Houston, TX, USA.

³⁴Division of Genetics and Genomics, Boston Children's Hospital, Harvard Medical School, Boston, MA, USA.

³⁵Division of Genome Sciences, National Human Genome Research Institute, Bethesda, MD, USA.

³⁶Department of Pediatrics, Division of Genetic Medicine, Seattle Children's Hospital, Seattle, WA, USA.

³⁷Department of Biomedical Data Science, Stanford University, Stanford, CA, USA.

Abstract

Rare diseases are collectively common, affecting approximately one in twenty individuals worldwide. In recent years, rapid progress has been made in rare disease diagnostics due to advances in DNA sequencing, development of new computational and experimental approaches to prioritize genes and genetic variants, and increased global exchange of clinical and genetic data. However, more than half of individuals suspected to have a rare disease lack a genetic diagnosis. The Genomics Research to Elucidate the Genetics of Rare Diseases (GREGoR) Consortium was initiated to study thousands of challenging rare disease cases and families and apply, standardize, and evaluate emerging genomics technologies and analytics to accelerate their adoption in clinical practice. Further, all data generated, currently representing ~7500 individuals from ~3000 families, is rapidly made available to researchers worldwide via the Genomic Data Science Analysis, Visualization, and Informatics Lab-space (AnVIL) to catalyze global efforts to develop approaches for genetic diagnoses in rare diseases (<https://gregorconsortium.org/data>). The majority of these families have undergone prior clinical genetic testing but remained unsolved, with most being exome-negative. Here, we describe the collaborative research framework, datasets, and discoveries comprising GREGoR that will provide foundational resources and substrates for the future of rare disease genomics.

Accelerating Diagnoses

The past decade has seen rapid progress in clinical genetics due to increased discovery of genes and variants involved in Mendelian diseases and ongoing advances in sequencing, variant analysis and data sharing¹⁻⁵. Despite this progress, most individuals who undergo clinical genetic testing for a suspected Mendelian condition remain undiagnosed⁶⁻⁹. For example, in the National Human Genome Research Institute (NHGRI) Centers for Mendelian Genomics, while over 3,800 genes were implicated in Mendelian disease, only about 11,000 out of over 28,000 families received a confirmed or potential molecular diagnosis. Thus, significant challenges remain to increase the molecular diagnostic yield and explain currently unsolved rare genetic disorders (**Box 1: Challenges in Diagnosing Rare Genetic Diseases**). In 2021, the NHGRI launched the GREGoR Consortium with five primary research sites and a data coordinating center to accelerate rare disease genetic research by harnessing the latest advances in sequencing including and especially genome sequencing and multi-omics; evaluating and prioritizing the use of functional genomics and novel computational strategies including recent advances in artificial intelligence; translating advances into routine clinical testing; advancing data sharing to foster a quorum of evidence for discovery; and collaborating with rare disease consortium worldwide to continue discovery and reporting of genetic etiologies for Mendelian diseases (**Fig. 1**).

Box 1: Challenges in Diagnosing Rare Genetic Diseases:

- A) The pathogenic variant(s) may be located in a gene yet to be implicated in disease. Until now, over 5,000 protein-coding genes^{10,11} have been implicated in at least one disease, but it is estimated that still 10,000+ disease gene relationships^{12,13} are undiscovered in just the remaining protein-coding genes.
- B) The pathogenic variant(s) may be located in the noncoding genome, where the mechanisms for how a variant manifests a clinical phenotype are not well understood leading to challenges in identifying candidates.
- C) The variant may be difficult to detect from solely short-read, exome or genome sequencing such as long repeats, inversions, and complex genomic rearrangements. The variant may be detectable but bioinformatic algorithms may struggle to call the variant correctly such as multi-nucleotide and mosaic variants. The variant may be detectable and called correctly but asserting its functionality or pathogenicity may require unavailable, orthogonal lines of evidence.
- D) More complex inheritance patterns such as multi-locus pathogenic variation, oligogenic, polygenic, variable expressivity, incomplete penetrance, imprinting, and/or mosaicism may also be confounding a diagnosis and necessitate a broader approach to understanding the mechanism of disease.
- E) A gene-disease relationship may be published or submitted to a genetic database but has yet to be reviewed and incorporated into clinical testing. This is compounded by the rapidly increasing number of Variants of Uncertain Significance (VUS).
- F) Many candidate variants and genes are n=1 regardless of the best data sharing practices.
- G) Because of the nature of novel discovery, there is not always a functional assay available to provide orthogonal evidence for or against a candidate variant or gene. Many times if a candidate does meet the inclusion criteria for an existing assay, the molecular phenotype measured in the assay may not match the potential mechanism of disease or fully recapitulate the pathophysiological impact, resulting in ambiguous results or an incorrect prediction of pathogenicity.
- H) Databases predominantly capture genetic information from individuals of European-like genetic ancestry potentially propagating biases in tools and reference data for variant classification for individuals of non-European-like genetic ancestry¹⁴.
- I) Newer genomic technologies may offer advantages over short-read DNA sequencing, but effective and widespread use of these technologies requires clear guidance and broad demonstration of efficacy.

EVALUATING EMERGING METHODS FOR ASCERTAINING RARE DISEASE DIAGNOSES

Squeezing the Exome

The most impactful approach to date for diagnosing rare diseases has been exome sequencing and periodic reanalysis of the protein coding sequences in the human genome¹⁵⁻¹⁹. GREGoR has led or contributed to 83 papers studying molecular diagnoses in 365 genes with more than a third being novel disease gene discoveries or phenotypic expansions²⁰⁻⁹⁹ (**Supplementary Table 1: Tracking GREGoR Papers With Molecular Diagnoses**) and provided a variety of automated pipelines for large cohort, exome and genome reanalysis, which include phenotypic data integration¹⁶⁻¹⁸. The success of reanalysis is largely driven by new disease gene discoveries and phenotypic expansions since the original analysis¹⁷, however new tools focused on reanalysis of well-known disease genes and loci have continued to yield diagnostic successes¹⁰⁰. For example, the inability to phase short-read exome (and also short-read genome) data can confound diagnoses of pathogenic compound heterozygous variants for recessive diseases. To overcome this, GREGoR contributed to a highly accurate method for inferring phase, and has calculated and released all pairwise phasing estimates and usage guidance for rare coding variants in exomes occurring in the same gene through the Genome Aggregation Database (gnomAD)^{101,102}.

Further, new computational approaches are increasingly able to identify structural variants and copy number variants from exomes. GREGoR researchers have developed tools to identify and implicate hundreds of pathogenic structural variant diagnoses from existing exomes that may have otherwise gone unsolved^{38,103-106}. Combined, GREGoR's efforts to continually extract diagnoses from existing exomes demonstrate the ongoing potential for continued innovation in genomic data reanalyses for rare diseases worldwide.

Short-Read Genome Sequencing

GREGoR has published a framework guiding usage of genomic technologies when genetic testing via panel or exome sequencing is inconclusive¹⁰⁷. The next step is typically short-read genome sequencing (srGS). srGS has already become widespread, with large-scale initiatives like the *All of Us* program¹⁰⁸ and UK Biobank¹⁰⁹ releasing srGS for nearly a million individuals. However, similar to exomes, the full diagnostic potential of srGS is still underutilized. For example, the SeqFirst-Neo program¹¹⁰ using rapid first-line, short-read genome sequencing based on broad eligibility criteria, obtained a precise genomic diagnosis in 50% of infants in an intervention group versus just 10% in the conventional care group. One year later, the intervention group had ninefold greater odds of diagnosis

compared to the control group and five times as many infants from underrepresented backgrounds received diagnoses. Ongoing GREGoR research aims to further extract diagnoses from srGS and demonstrate its increasing utility as a first-line clinical test.

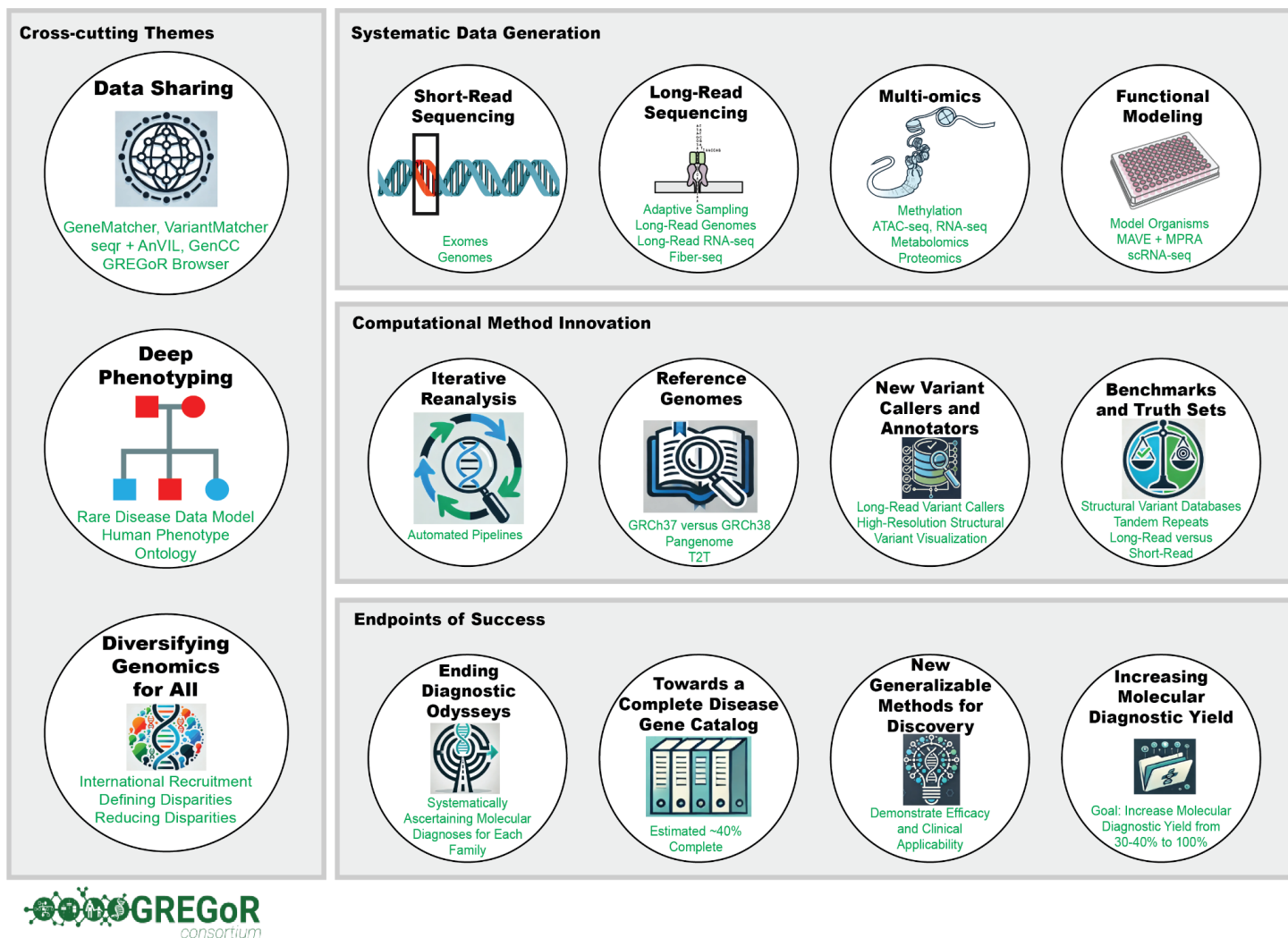


Figure 1 | Overview of GREGoR. Strategic framework of the GREGoR Consortium for accelerating genomics in rare disease research, highlighting cross-cutting themes, systematic data generation, computational innovations, and endpoints of success.

However, most molecular diagnoses deduced by srGS are found in protein-coding genes, suggesting they could potentially be detected by exome sequencing. To evaluate the relative utility of srGS, a large-scale study of 822 families by GREGoR researchers reported that of 218 patients who received a diagnosis via srGS, 72% of variants should have been detectable by exome sequencing¹¹¹. The remaining 28% of cases were explained by variants not readily accessible on exomes such as tandem repeat expansions, deep intronic variants, structural variants, and variants in difficult-to-sequence coding regions. Overall, srGS resulted in a >8% increase in the diagnostic yield compared to just exome sequencing and underscored the growing support for using srGS as a first-tier test.

GREGoR is also developing new visualization tools for structural and copy number variants¹⁰⁶ by repurposing read depth data from srGS to mimic SNP arrays to achieve resolution as low as 1 kb - beyond the 5 kb limit of the current standard using array comparative genomic hybridization. Thus, srGS could potentially serve as a cost-effective, first-line, unifying assay by simultaneously replacing both arrays and exomes and enable more accurate, nucleotide-resolution breakpoints of structural variants, which have historically been critical in uncovering mechanisms of genomic rearrangements. Most breakpoints including published structural variants lack validation at nucleotide resolution which is relevant for genomic assembly and hypothesis-driven inferences of SV impact on gene expression. In this same vein, GREGoR investigators are showing that local sequences surrounding candidate and pathogenic variants can offer insights into secondary structure mutagenesis and other mechanisms of genomic disorders^{70,71,112}. GREGoR's existing and ongoing work to uncover diagnoses from srGS emphasizes the substantial, yet underutilized, potential of both primary analysis and reanalysis of srGS for rare disease discoveries.

Long-Read Sequencing

Long-read sequencing has recently ushered in new diagnostic opportunities in rare diseases. GREGoR and others have demonstrated that use of targeted long-read sequencing can reveal variants, particularly structural variants spanning repetitive sequences, in both known and novel disease genes that are missed or difficult to detect by short-read sequencing^{40,97,113-115}. With

emerging long-read technologies that allow targeting of specific genomic regions such as adaptive sampling, targeted sequencing panels can evolve beyond coding regions to include UTRs, promoters, intronic, intergenic, and large expanses of noncoding regions around known disease genes. For example, GREGoR in collaboration with Twist developed the Twist Alliance Dark Genes Panel to produce phased variants across 389 medically-relevant and complex autosomal genes, where short-read sequencing tends to fail¹¹⁶. Not only were novel pathogenic variants discovered, but an annotated resource was also created to address gaps in current databases for these genes.

One of the major challenges with use of long-read sequencing in rare diseases remains the absence of control datasets for filtering and prioritizing variants. GREGoR is using long-read genome sequencing from individuals with diverse genetic ancestry with an aim to create benchmarks and a database of structural variants for filtering and prioritization, and to catalog structural variants in genes that are difficult to sequence using short-read technology^{117–120}. Specifically, GREGoR researchers started the 1000 Genomes Project ONT Sequencing Consortium to generate a database of structural variants derived from long-read sequencing for filtering and prioritization of structural variants in unsolved individuals and recently released the first 100 samples of long-read data from diverse populations¹²¹. GREGoR is also establishing a baseline catalog of complex structural variants using optical genome mapping. Moreover, GREGoR is developing tools for improved and useful annotation for kilobase and megabase scale variants especially using long-read sequencing with a focus on mosaic structural variants^{122,123}. Additionally, innovative computational tools in long-read sequencing variant calling and analysis^{123,124} are being developed including *de novo* variant callers and long-read pipelines for mitochondrial variant calling.

As long-read sequencing becomes more commonly used, a key focus for GREGoR is comparing the molecular diagnostic yield and cost-effectiveness of short-read versus long-read sequencing for rare disease cohorts. In multiple studies, long-read genome sequencing uncovered novel candidate variants and genes that were missed by exome or short-read genome sequencing, including *de novo*, compound heterozygous, structural, and epigenetic variants^{69,115,125}. In comparison to short-read sequencing, long-read sequencing offers advantages such as better phasing, improved understanding of haplotype blocks, and methylation analysis. GREGoR researchers have been leveraging these advantages by developing tools¹²⁶ to utilize methylation data and investigating the diagnostic yield improvements from genome-wide DNA methylation arrays in relation to long-read sequencing¹²⁷.

In addition to methylation analyses, usage of long-read sequencing to offer multi-omic insights beyond traditional DNA sequencing are increasingly being applied to find and understand molecular diagnoses as well as mechanisms in rare diseases. One such technology is Fiber-seq, which uses long-read sequencing to simultaneously evaluate primary DNA sequence with a nucleotide-resolution view of surrounding chromatin and epigenetic architecture^{128,129}. GREGoR is developing unifying assays looking to simultaneously assay the genome, methylome, epigenome, and transcriptome to identify and understand mechanisms of rare diseases¹³⁰. Such unifying assays help explain previously elusive variants⁹⁵ that may have been visible on exome or short-read genome sequencing but may not have been nominated as candidate variants. The comparison and contrast of multiple layers of -omics data in a single, unifying assay mechanistically implicates the pathogenicity of these candidates, especially for noncoding variants.

Multi-omics

Complementing advances in DNA sequencing, rapid advances in development of -omics assays continue to provide insights into genome function. In rare disease diagnosis, methylome and transcriptome data have broadly demonstrated their utility by identification of outlier events in methylation¹³¹, splicing, or gene expression implicating pathogenic variants^{132–134}. However, there is no standardization of clinical and research transcriptome sequencing despite its promise as a primary diagnostic tool¹³⁵. GREGoR has focused on complementing hundreds of cases with methylation and transcriptome data to facilitate development of standards and new computational methods. For example, recent activities in GREGoR have demonstrated how combined transcriptome and long-read genome analyses can aid in prioritizing structural variants when allele frequency information is limited¹³⁶.

An ongoing focus of GREGoR has been creating data that enables evaluation and prioritization of multi-omic assays for rare disease diagnosis. Currently, usage of multi-omics is predominantly limited to research and limited information exists to suggest which post-genome, -omics assay would yield the most useful information. To address this challenge, GREGoR has been generating a squared-off matrix for a subset of families, for whom long-read genome, methylation, chromatin-accessibility, transcriptome, proteome and metabolome data are being collected. Complementing these data, GREGoR has been supporting development and integration of reference -omics data from the Common Fund Data Ecosystem to advance outlier detection for various -omics assays by integrating larger control datasets. Multiple efforts in GREGoR are also facilitating more routine use of these data such as updates to the *seqr* platform to expand intake of multi-omics data types to enable routine linking of outliers in these data to underlying genomic variation¹³⁷.

Functional Modeling

Functional assays have been critical to identify, screen, and/or validate candidate genes and variants to confirm novel genotype-phenotype relationships and phenotypic expansions and to uncover underlying genomic mechanisms of disease. Of the 83 papers led or facilitated by GREGoR resulting in novel disease gene discoveries or phenotypic expansions, 44 were supported by orthogonal functional experiments to confirm these findings ([Supplementary Table 1: Tracking GREGoR Papers With Molecular Diagnoses](#)). Much of the functional work in the rare disease field has historically relied on classic model organisms like zebrafish, yeast, fruit flies, and mice. Model organism research continues to play a significant role in GREGoR's collaborative efforts with the broader genomics and scientific community, advancing the understanding of disease genes and their variations¹³⁸.

An emerging area of technologies for rare disease research is in single cell studies. For example, GREGoR researchers have been able to make high-resolution maps of fetal hematopoiesis to understand how Trisomy 21 predisposes to hematological malignancies¹³⁹; understand regulatory programs contributing to hair and skin diseases¹⁴⁰; and discover candidate causal variants for Alzheimer and Parkinson disease¹⁴¹. Additionally, GREGoR researchers are developing single cell technologies to simultaneously profile chromatin accessibility, transcriptomics, and nuclear protein abundance¹⁴² as well as applying improvements on single-cell, whole genome amplification methods to understand somatic copy number variations in healthy and diseased brain tissue¹⁴³.

In addition to single cell and model organism studies, new functional experiments and high-throughput functional genomics are progressively integrated into GREGoR's research through cross-consortium collaborations with the NHGRI IGVF (Impact of Genomic

Variation on Function) consortium¹⁴⁴. Technologies such as Multiplexed Assays of Variant Effect¹⁴⁵, Massively Parallel Reporter Assays¹⁴⁶, mini-gene splicing assays¹⁴⁷, and high-throughput imaging for cellular mislocalization¹⁴⁸ are providing functional data to validate or invalidate candidate genes and variants and give mechanistic understanding to variant penetrance and pleiotropy^{149,150}. For example, GREGoR has exported over 500 candidate genes and variants to IGVF, and more than 100 have been selected for further functional evaluation including in the IGVF's Perturb-seq experimental plans. Further, GREGoR investigators are generating hundreds of isogenic human induced pluripotent stem cell (hiPSC)-derived neural stem cells and glutamatergic induced neurons with CRISPR-engineered, systematic structural variant deletions of local topologically associated domains and chromatin loops^{151,152}. These models are eligible to be shared as a GREGoR resource, and in conjunction with transcriptomic and single cell profiles from mice, are being used to study the complex and context-dependent impacts of structural variation on neurodevelopment. Finally, GREGoR is importing predictions of noncoding variant functions from IGVF to improve the understanding of noncoding mechanisms involved in rare diseases¹⁵³. At the same time, GREGoR researchers are advancing the creation and use of developmental cell atlases to develop deep learning models trained on chromatin accessibility and gene expression data representing diverse adult, fetal and developmental contexts to enable the identification of context-specific regulatory effects of rare and *de novo* noncoding variants.

REFRAMING RARE DISEASE ANALYSIS

Reference Genomes

Adoption of new reference genomes has lagged in clinical settings. Despite publication of GRCh38 over a decade ago, many clinical labs till today still use GRCh37¹⁵⁴. Part of the entrenchment of GRCh37 was the lag in necessary infrastructure development to support allele frequencies, *in silico* scores, bioinformatic tools, and clinical databases on GRCh38. Thus, till today, the majority of known clinical disease genes and phenotypic expansions were initially discovered using GRCh37. As a result, research groups, including those in GREGoR, have been studying the differences between GRCh37 and GRCh38 and more references in variant calling and downstream analyses. At the exome level, it has been shown that the reference genome alone impacts variant calling in ~1% of the exome, with 206 genes enriched in discordant calls, including 8 known disease genes¹⁵⁵. These discrepancies were more pronounced at the RNA level, with research in GREGoR highlighting that 1,492 genes demonstrate reference-dependent quantification, 3,377 genes exhibit reference-exclusive expression, affecting 512 known disease genes¹⁵⁶. GREGoR investigators have also focused on fixing the GRCh38 reference¹⁵⁷, benchmarking medically relevant genes for both GRCh37 and GRCh38¹⁵⁸, and resolving pathogenic inversions in reference genome gaps using the telomere-to-telomere (T2T) reference¹²⁵.

An important question for the field focuses on whether there will be development of flexible pipelines and tools capable of using the newest references, such as T2T and the pangenome¹⁵⁹. GREGoR in collaboration with Illumina has benchmarked the DRAGEN pipeline which uses graph-based alignment among many other novel features for variant calling in short-read genome sequencing¹⁶⁰. Looking ahead, GREGoR is collaborating with the Human Pangenome Reference Consortium (HPRC) through methods development like the Pangenome Research Tool Kit¹⁶¹ to demonstrate accurate variant calling in regions of the genome that were previously too complex for accurate variant calling. Further GREGoR investigators are utilizing pangenome approaches to understand complex tandem repeats in known disease genes¹⁶² and exploring the infrastructure necessary for widespread adoption of the pangenome in clinical settings.

Deep Phenotyping

Studying phenotypic heterogeneity in the context of genetic heterogeneity is critical to solving unsolved Mendelian disease. Assignment of Human Phenotype Ontology (HPO) terms¹⁶³ is a mandatory requirement for GREGoR data collection to allow end users to link all possible genotypes to all possible phenotypes. Additionally, GREGoR is developing novel algorithms based on the directed, acyclic HPO graph to resolve blended phenotypes resulting from multilocus pathogenic variation^{8,56}, implicate genetic heterogeneity as resulting in similar phenotypic manifestations^{49,71}, and elucidate gene- and variant-driven phenotypic heterogeneity in rare diseases that have demonstrated genetic heterogeneity^{80,164}. Upstream of phenotypic analysis pipelines, GREGoR is building large language models for optimal phenotypic extraction from electronic health records¹⁶⁵. Understanding phenotypic complexity is particularly critical in tackling the rarest and most challenging molecular diagnoses (**Box 2: The Rarest and Hardest Molecular Diagnoses**), where genetic heterogeneity and diverse mechanisms of pathogenicity require nuanced interpretation and integration of multiple technologies.

Box 2: The Rarest and Hardest Molecular Diagnoses:

Over time, consistent themes have emerged in the life cycle of disease gene discovery. Typically, the first and easiest candidate variants implicated are *de novo* and/or predicted loss-of-function (pLOF) variants that segregate with the phenotype in a pedigree. These two variant classes have logical frameworks supporting their putative mechanisms. A *de novo* variant, found in an affected proband but not in the unaffected parents, significantly increases the probability of being causative, especially when other *de novo* cases show the same clinical phenotypes^{166,167}. pLOF variants such as nonsense single nucleotide variants, frameshifting insertions or deletions, and splice-site altering variants imply a null effect, because many of these pLOF variants are expected to be caught by the mRNA surveillance mechanism known as nonsense-mediated decay (NMD)^{168,169}. NMD is a highly sensitive mechanism present in all tissues and destroys faulty transcripts with premature stops with high efficiency and fidelity^{170,171}. This same mechanistic reliability explains why pLOF variants are often the first variant type to be implicated in a novel gene-to-disease discovery, because clinical geneticists can reliably infer that the presence of a pLOF variant will putatively lead to destruction of the faulty RNA transcript and no protein production which is in alignment with a potential pathogenic mechanism of loss of function. Subsequently, other single nucleotide and structural variants in the same gene or region are often implicated after a quorum of cases is established although different types or locations of variation in the same gene can result in distinct clinical phenotypes. Further, once one gene has been implicated it serves as a seed for other genes in the same protein complex¹⁷² or protein pathway⁷⁰ or gene family^{173,174} to be implicated in the same or similar clinical phenotypes.

Currently, Online Mendelian Inheritance in Man (OMIM) and the Gene Curation Coalition (GenCC) have documented over 5000 genes as being implicated in at least one Mendelian condition^{10,11}. Most of these discoveries were achieved through sequencing and interpretation of primary DNA variation in the context of a proband's phenotypes without requiring additional -omics or integrative analyses. While many thousands more disease genes can still be discovered with these same established gene discovery principles, GREGoR is pursuing cases unsolved by standard clinical genetic testing hypothesized to have among the rarest and most difficult-to-detect or interpret pathogenic variation, which many times requires integration of multiple -omics to 1) discover a candidate variant refractory to traditional sequencing methods; 2) provide proper context to interpret a variant that may have been seen in the primary DNA sequence but there was not enough understanding to nominate the variant as a candidate; or 3) provide orthogonal validation of a candidate discovered in the primary DNA sequence but for which the interpretation was speculative at best. Below we discuss 12 of the rarest and hardest molecular diagnoses pursued by GREGoR investigators:

- 1. Noncoding variants** occur in genomic regions that do not code for proteins, such as noncoding RNAs, promoters, enhancers, and untranslated regions. Despite not altering protein sequences, these variants can disrupt mechanisms such as gene regulation, splicing, expression, or RNA stability, playing a role in development of rare diseases. Genome sequencing is crucial for detecting variants in noncoding regions, but techniques like ChIP-seq, ATAC-seq, RNA-seq, Fiber-seq and Massively Parallel Reporter Assays can help identify and provide mechanistic explanation for potential pathogenic noncoding variation. For example, noncoding variants such as deep intronic variants may be involved in alterations in splicing, where exons or introns may be incorrectly skipped or included during mRNA processing leading to changes in the final transcript potentially leading to clinical phenotypes¹⁷⁵. While these variants are often visible in primary DNA sequencing, their interpretation is typically speculative without orthogonal validation such as RT-PCR, RNA sequencing, or mini-gene splicing assays to validate whether or not aberrant splicing occurred by understanding the sequence of the processed mRNA. Similarly, long-read RNA-sequencing is particularly useful for detecting full-length transcripts and interpreting complex splicing patterns. Additionally, splicing-specific variant effect predictors¹⁷⁶ and variant effect predictors that tabulate scores for all the >9 billion single nucleotide variants in both coding and noncoding regions can be deployed to identify potentially pathogenic noncoding variants. Work by GREGoR and others have shown that noncoding variants can be the 'missing variant' in *trans* with a pathogenic coding variant for recessive rare diseases^{177,178}. GREGoR is pursuing generalizable methods for understanding noncoding variant effects at scale^{44,93}.
- 2. Noncoding genes** produce molecules that perform potential regulatory, structural, or catalytic roles rather than encode proteins. These include rRNA, tRNA, miRNA, lncRNA, snRNA, and more, and perturbations in noncoding genes may cause rare genetic diseases. A flagship example is perturbations in *RNU4-2* in which cases from multiple GREGoR sites contributed to rapid discovery and progress to publication establishing perturbations in *RNU4-2* as one of the most commonly mutated causes of neurodevelopmental disorders^{75,179}. Even though *RNU4-2* is a noncoding RNA, its discovery timeline is analogous to the typical gene discovery life cycle. The main cases found were *de novo* insertions at the same site in a very large cohort with phenotype-matching cases. *RNU4-2* has now served as a seed for more *RNU4* minor spliceosome genes being implicated in neurodevelopmental disorders^{180,181}. Another example from GREGoR is *CHASERR*, a long noncoding RNA adjacent to *CHD2* which was implicated in developmental and epileptic encephalopathy⁹³. The *CHASERR* discovery was a strong collaboration with the father of the initial proband serving as a coauthor on the manuscript, highlighting the power of patient partnerships in accelerating rare disease genomics. Further, using Fiber-seq, GREGoR investigators have identified *STRTS*, an intergenic locus implicated in congenital hypothyroidism⁹⁵. These findings illustrate the diagnostic potential of noncoding regions of the genome, which are not systematically included in standard variant analysis workflows. With thousands of noncoding transcripts still poorly understood, GREGoR continues to explore this untapped reservoir of genomic information, paving the way for novel disease-gene discoveries in the noncoding space.
- 3. Multilocus Pathogenic Variation (MPV)** refers to the presence of multiple, independently pathogenic variants in multiple genes or loci that collectively contribute to an individual's clinical manifestation¹⁸²⁻¹⁸⁴. These variations can interact in complex ways, leading to compounded effects that may influence disease severity, onset, or progression, often complicating diagnosis and treatment¹⁸⁵. GREGoR researchers have shown that as many as 5% of individuals where molecular diagnoses are ascertained have MPV⁸, with this rate being even higher in rare disease families with parental consanguinity^{38,184}. While individual variants comprising MPV are typically routinely diagnosed from just DNA sequencing, the interpretation of MPV typically requires deeper phenotypic analyses and represents a much broader area of gene dosage models for disease causing variation¹⁸⁶.
- 4. Oligogenic molecular diagnoses** involve the contribution of variants in two or more genes or loci that collectively lead to a clinical phenotype, a concept distinct from traditional monogenic inheritance patterns or MPV. These cases often present diagnostic challenges because the individual variants may not cause disease independently but act in concert to disrupt pathways or biological networks. GREGoR investigators have shown especially for congenital heart disease the importance of digenic and oligogenic mechanisms^{94,182,183}. Currently, we can identify oligogenic variants through DNA sequencing, but the full scope of oligogenic molecular diagnoses may extend across multiple layers of multi-omics. For example, combinations of transcriptomic, proteomic, or metabolomic alterations may converge with DNA variation to create synergistic effects that contribute to disease, representing a new frontier for rare disease genomics.
- 5. The Multiple *de novo* Copy Number Variant (MdnCNV)** phenotype is a form of multilocus pathogenic variation where four or more independent, constitutional *de novo* copy number variants arise in the same person within one generation^{187,188}. GREGoR researchers have shown that this ultra-rare phenotype occurs in roughly 1 in over 12,000 individuals referred for genome-wide chromosomal microarray analysis. MdnCNV typically requires integration of multiple technologies as well as use of quantitative phenotyping analysis to fully characterize and understand genomic and clinical impact, including but not limited to arrays, short-read sequencing, and long-read sequencing.
- 6. Complex Genomic Rearrangements (CGRs)** are kilobase to megabase-scale structural changes in the genome that involve multiple breakpoints, rearrangements, and/or integration of novel sequences in *cis* that result in duplications, deletions, inversions,

and translocations, often affecting gene function and regulation. While CGRs have been catalogued^{189–192} and shown to be very abundant across diverse populations, the accurate detection and assembly of CGRs, especially their breakpoints, typically requires more than just short-read DNA sequencing. Long-read sequencing (including adaptive sampling and ultra long-reads), optical genome mapping, chromosomal microarray analysis, and linked-read sequencing can help provide nucleotide resolution for CGRs to understand their mechanisms of formation to better understand how they precipitate genomic disorders and contribute to clinical variability and disease severity^{193–195}.

7. **Tandem repeats** are sequences where a nucleotide motif is repeated consecutively a varying number of times^{196–199}. These repetitive regions can be unstable, leading to expansions or contractions, which are associated with several genetic disorders. Due to the potential for multi-mapped short reads, only shorter tandem repeats have typically been well detected from short-read sequencing technologies. However, long-read sequencing technologies are particularly effective at completely spanning short and longer tandem repeats, enabling accurate determination of repeat length and structure. Additionally, specialized bioinformatic tools^{162,200,201} designed for repeat analysis can help in accurately calling tandem repeats and identifying pathogenic expansions or contractions. GREGoR is collating large databases and truth sets of tandem repeats across diverse populations^{202–204} to enable more systematic integration of tandem repeat analysis into sequencing pipelines.
8. **Mosaic Variation** originates from post-zygotic mutations and is considered germline if confined to the germ cells or somatic if acquired during or after the first mitotic divisions. Mosaic variation can lead to variations in phenotype, depending on the proportion and distribution of the mutant cells across tissues. Detecting mosaic variation often requires more than standard DNA sequencing due to the low variant allele fraction of mosaic alleles. Reliable detection of mosaic variation typically requires high read-depth sequencing and orthogonal techniques such as digital droplet PCR for validation and bioinformatic discrimination¹²³ between mosaic variants and sequencing artifacts. GREGoR is collaborating with the SMaHT (Somatic Mosaicism across Human Tissues) consortium to understand pathogenic mosaic variation at scale.
9. **Multi-nucleotide Variants (MNVs)** are two or more variants within the same codon on the same haplotype^{205–207}. Accurately identifying and interpreting MNVs is typically a challenge for variant callers and annotation pipelines that may incorrectly interpret the single MNV as multiple independent variants. However, it has been shown that there are over 50 variants per person affected by MNVs²⁰⁸. Incorrect interpretation of MNVs can alter the interpretation of clinically pathogenic variation such as nonsense single nucleotide variants that do not lead to a premature stop codon introduction. Many variant effect predictors and Multiplex Assays of Variant Effects systematically score every possible MNV in a target locus. GREGoR is evaluating the utility of these data towards VUS reclassification and novel disease gene discovery.
10. **NMD Escaping Variants** introduce premature stop codons in mRNA that evade NMD, producing truncated proteins of unknown function that may or may not manifest clinical disease. These variants are visible on primary DNA sequencing and because extensive work has been done to determine the rules of NMD escape^{209,210}, many of these variants are already implicated in clinical disease and can be speculated on from just primary DNA sequencing^{171,211}. However, newer approaches such as long read RNA-seq and proteomics are opening new doors into the investigation of NMD escape alleles and their downstream mechanisms.
11. **Incompletely Penetrant Variants** refer to pathogenic changes in DNA that do not always result in observable clinical disease, even in individuals carrying the variant. This variability can complicate interpretation, especially in family studies where carriers appear unaffected^{68,71,212}. These variants often challenge diagnostic workflows because traditional penetrance assumptions do not hold, necessitating integration of orthogonal lines of evidence such as animal models^{71,212} and epigenetics⁸⁹ to understand the variant pathogenicity. These functional efforts are particularly critical for diseases where penetrance may be age-dependent, sex-influenced, or modified by external factors⁸⁹. Further, in a case-by-case analysis across all of gnomAD v4, >95% of incompletely penetrant pLOF variants found in severe, early-onset, highly penetrant haploinsufficient disease had explainable causes such as a downstream frame-restoring variant, predicted re-initiation by a downstream methionine, a MNV changing the interpretation of a nonsense variant to a missense or synonymous variant, or the location of the pLOF variant being in an NMD Escape region²¹³.
12. **Variants of Uncertain Significance (VUS)** are most often reported from testing genes that already have been established in disease pathogenesis²¹⁴, though by definition, all variants found in candidate genes with insufficient evidence for disease implication are also VUS. VUS are accumulating rapidly over time as testing volume expands. In fact, VUS are more often reported during panel testing compared to exome or genome sequencing due to professional practices²¹⁵ and are disproportionately called in individuals from non-European-like genetic ancestry¹⁴. Thus, while transitioning to consistent first-line clinical usage of exome or genome sequencing coupled with phenotypic analyses will decrease the rate of clinically reported VUS, given that most causal variants are only found in a single individual, integration of functional modeling is often required to reclassify VUS. Multiplexed Assays of Variant Effects are high-throughput experiments orthogonal to the clinical sequencing pipeline that produce functional scores for all variant effects in a target locus and when their evidence strength is clinically calibrated and incorporated with other lines of evidence, demonstrate significant promise in massive VUS reclassification^{14,216}.

DIVERSIFYING RARE DISEASE GENOMICS AND DIAGNOSTICS

Participation in rare disease research can be influenced by numerous factors, including but not limited to institutional, socioeconomic, geographic, linguistic, cultural, educational, and insurance factors^{217–220}. GREGoR sites, many of which are in urban centers, have implemented procedures for online enrollment, chatbots²²¹, remote consent, offsite sample collection (including mobile phlebotomy in rural areas), and translated materials in multiple languages to improve accessibility. Attention to these details has allowed GREGoR to foster many international collaborations with local scientists to sequence and make available sequencing data from thousands of individuals of non-European-like genetic ancestry, actively seeking to address the disparities in genomic data availability across Middle Eastern, North African, Southeast Asian, South American, and other underrepresented groups such as African-American and Hispanic peoples.

Diversifying participants in genomics research via recruitment of participants from underrepresented populations is just one approach to fostering equity. GREGoR is pursuing orthogonal approaches to increase access to a genetic diagnosis by pursuing improved variant calling methods, applying multiplexed functional assays to improve interpretation of VUS that are enriched in underrepresented populations, and testing technologies and workflows that reduce barriers to equitable access to a genetic diagnosis. For example, GREGoR is collaborating with the Human Pangenome Reference Consortium to reduce bias and improve variant calling accuracy across all populations via usage of the pangenome. GREGoR has analyzed population biobank data to show a higher prevalence of Variants of Uncertain Significance and fewer Pathogenic or Likely Pathogenic classifications in individuals of non-European-like genetic ancestry¹⁴. These disparities were alleviated by using high throughput, multiplexed functional experiments to test every possible single variant in genes of interest to resolve VUS disparities between populations. However, this study demonstrated that allele frequency and variant effect predictors contribute to the inequitable classification of variants and more work to prevent bias in clinical variant classification is an important future priority for GREGoR. Finally, the SeqFirst program¹¹⁰ shows that using simple criteria to assess eligibility for rapid short-read genome sequencing significantly increases the proportion of non-White and Black infants who receive a precise genetic diagnosis. GREGoR and SeqFirst are conducting a comparison of long-read versus short-read sequencing in the SeqFirst cohort to better understand the relative value of these technologies within populations in this cohort.

ACCELERATING DATA SHARING

Data sharing is critical to the advancement of rare disease genomics and diagnoses^{222,223}. GREGoR is committed to rapid release of genomic and phenotype data to the larger research community within AnVIL and making these data FAIR (Findable, Accessible, Interoperable, and Reusable)²²⁴ and machine-readable. At the time of AnVIL submission and prior to analysis, deep phenotyping (Fig. 2A), potentially multiple orthogonal lines of -omics data (Fig. 2B, 2C), and pedigree data (Fig. 2D) are made available via the GREGoR data model for broader dissemination to the research community. Currently, DNA data on approximately 7400 individuals from over 3000 families is available with transcriptome data available for over 500 individuals and nearly 200 participants with both exome and short-read genome data (dbGaP:phs003047; Fig. 2B). Within the next year, planned releases include additional short and long read genomes, short and long read RNA-seq, Fiber-seq, ATAC-seq, metabolomics and proteomics. GREGoR has begun to assess cases (Fig. 2E) and has identified candidate discoveries and molecular diagnoses for over 400 families which are added in subsequent AnVIL submissions over time. GREGoR hopes that solved cases can be used by the community as positive controls for benchmarking tools and that the current dataset of unsolved cases (many of which are exome negative) holds tremendous potential for discovery by collaborative community efforts.

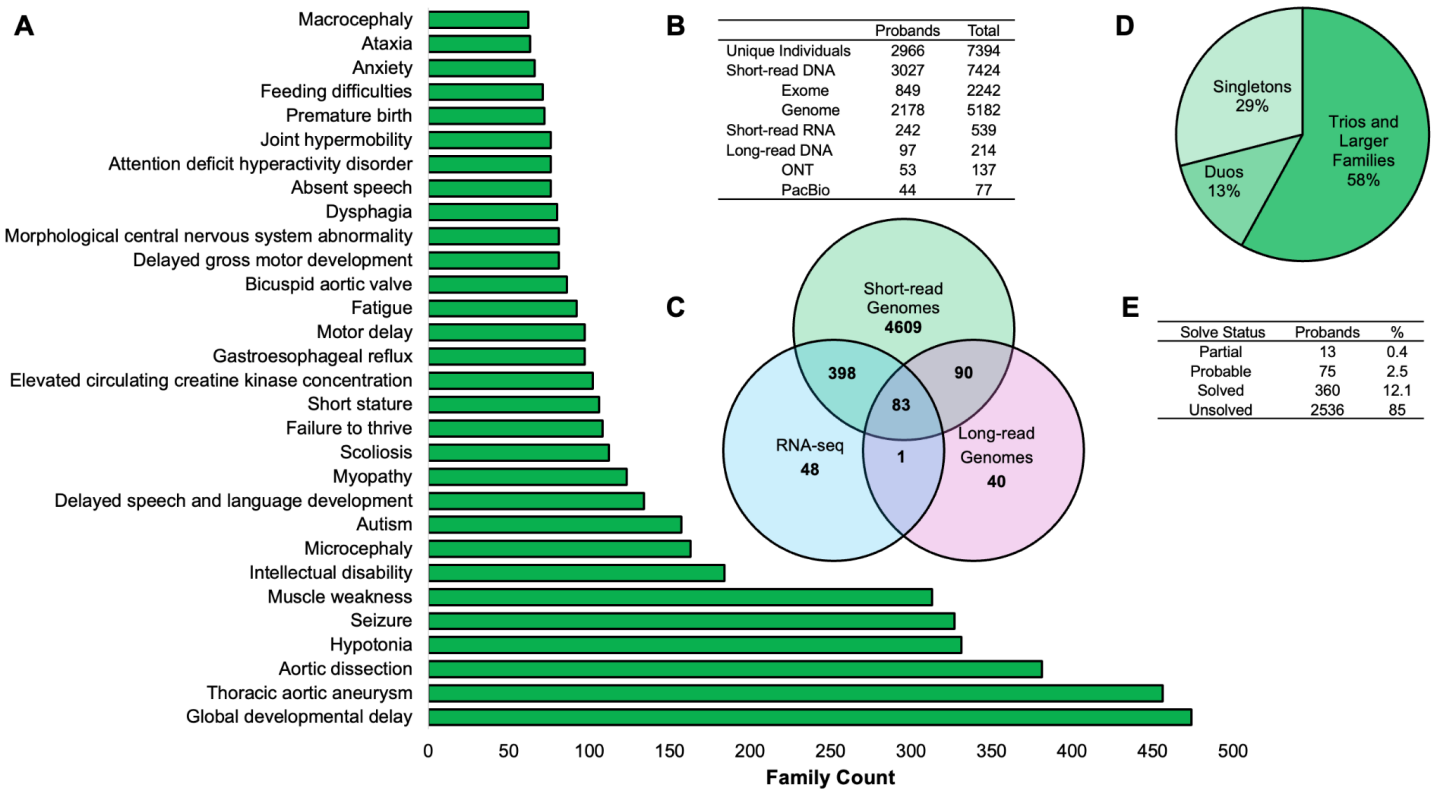


Figure 2 | Overview of Publicly Released GREGoR Data. Summary of second public data release (dbGaP:phs003047). (A) Distribution of top 30 phenotypes in GREGoR based on Human Phenotype Ontology descriptions. (B) Table of numbers for probands and total individuals for each sequencing modality. (C) Venn diagram depicting overlap across short-read genomes, RNA-seq, and long-read genomes in data generation. (D) Family structures comprising the overall cohort from a total of n=3059 families. (E) Summary of current solved cases. Data is shared prior to analysis but even the current diagnostic outcomes underscore the challenges and opportunities in resolving rare disease cases that are previously exome negative.

Many commercial and academic clinical labs have undiagnosed cases potentially explained by novel gene discoveries or phenotypic expansions. The GREGoR Consortium has actively engaged with clinical labs to study effective strategies for exome and genome analysis without overburdening variant analysts²²⁵. GREGoR has released recommendations for clinical labs to report variants in novel candidate genes and support follow-up investigations, enabling broad discoveries and patient diagnoses²²⁶.

To enable exchange of data and facilitate collaboration, all GREGoR candidate genes are shared to Matchmaker Exchange^{227,228} through either GeneMatcher²¹², *seqr*¹³⁰ or MyGene2²²⁹. Notably, for the 83 GREGoR publications ([Supplementary Table 1: Tracking GREGoR Papers With Molecular Diagnoses](#)) involving novel disease genes or phenotypic expansions, almost every project has been influenced by findings from connections made across or within nodes of the Matchmaker Exchange. Novel candidate genes and phenotypic expansions are also curated for validity and publicly shared to the Gene Curation Coalition database to accelerate access to early evidence of novel gene-disease relationships and aid in standardized clinical diagnostics and research¹¹. Analogously, candidate variants and molecular diagnoses are deposited in ClinVar^{230,231}. Also GREGoR is leveraging federated variant-level matchmaking through tools such as VariantMatcher^{232,233}, which allows queries of variant data across different genomic datasets, even when a variant or gene has not been recognized as a candidate with the hope of accelerating disease-causing variant-level discovery within and beyond the exome.

GREGoR has developed a novel data model that emphasizes the essentials of rare disease research such as the importance of data accessibility, consent consistency and transparency, and the usage of many accepted ontologies and common standards. Every variant in the GREGoR joint callset is machine-readable with both a unique GA4GH Variant Representation Specification (VRS) ID²³⁴ and ClinGen allele ID²³⁵, and the data model accommodates variants and output files from a wide variety of genomic, multi-omic, phenotypic, and molecular data types, and is modularly designed to support the integration of future data types. GREGoR implements this model within AnVIL workspaces for data submission and validation to enable (1) streamlined genomic and phenotypic data deposition; (2) scalable, semi-automated quality control of molecular, phenotypic, and variant annotation data; and (3) expedited, controlled-access release to the broader research community.

Inside AnVIL, the GREGoR data is also queryable via *seqr*, which integrates variant filtration, annotation, and causal variant identification¹³⁷. Outside AnVIL, GREGoR has developed a public variant browser²³⁶, which already includes >95 million variants. Not only is the number of families predicted to triple by the end of GREGoR, but de-identified phenotypic data are also being added to the public browser, allowing researchers to easily explore putative genotype-phenotype relationships in rare disease families.

CONCLUSION

GREGoR aims to advance state-of-the-art approaches to determine molecular diagnoses for individuals with unsolved rare diseases. Central to this approach is the creation of a broadly accessible and information-rich genomics data resource derived from individuals and families with rare diseases for whom prior standard-of-care testing such as exome sequencing has been non-diagnostic. This resource is defined by a future-forward data model and infrastructure that incorporates genomic and other -omic datasets generated through emerging technologies, family structure data, and rich phenotypic data and currently is supporting data for over 3500 families, many of whom remain unsolved.

Several opportunities remain to advance rare disease research and accelerate diagnoses. Foremost, the effort to complete a comprehensive catalog of genes underlying Mendelian conditions remains far from finished. While thousands of Mendelian conditions have been described, a significant proportion of these conditions still lack a known genetic cause, leaving substantial gaps in our understanding. Further, even when all genes for Mendelian conditions are identified, causal variant discovery will remain far from saturated. This is particularly true for missense variants, which often require functional validation, and for noncoding variants, where the regulatory mechanisms are complex and poorly characterized. Far from being a task of "wrapping up the edges," these challenges represent a vast forefront in genomic research, demanding both innovative methodologies and sustained collaboration to make meaningful progress. Alongside these challenges, the advancement of functional genomic assays has required the vetting of these approaches at scale in individuals of diverse genetic ancestries and diverse rare disease phenotypes. Such efforts are critical to establishing the standards of when and how to use a specific approach and will guide expectations on their relative yields at scale and their adoption in clinical practice. Lastly, there is a palpable need to translate scientific discoveries into curation practices that align with formal clinical standards. To address these gaps, GREGoR provides data and infrastructure that will catalyze the development and implementation of new approaches to advance genomics in rare disease by the broader community.

References

1. Baxter, S. M. *et al.* Centers for Mendelian Genomics: A decade of facilitating gene discovery. *Genet. Med. Off. J. Am. Coll. Med. Genet.* **24**, 784–797 (2022).
2. Taylor, J. C. *et al.* Factors influencing success of clinical genome sequencing across a broad spectrum of disorders. *Nat. Genet.* **47**, 717–726 (2015).
3. Wright, C. F. *et al.* Genomic Diagnosis of Rare Pediatric Disease in the United Kingdom and Ireland. *N. Engl. J. Med.* **388**, 1559–1571 (2023).
4. Rimmer, A. *et al.* Integrating mapping-, assembly- and haplotype-based approaches for calling variants in clinical sequencing applications. *Nat. Genet.* **46**, 912–918 (2014).
5. Monaco, L. *et al.* Research on rare diseases: ten years of progress and challenges at IRDiRC. *Nat. Rev. Drug Discov.* **21**, 319–320 (2022).
6. Yang, Y. *et al.* Molecular findings among patients referred for clinical whole-exome sequencing. *JAMA* **312**, 1870–1879 (2014).
7. Farnaes, L. *et al.* Rapid whole-genome sequencing decreases infant morbidity and cost of hospitalization. *NPJ Genomic Med.* **3**, 10 (2018).
8. Posey, J. E. *et al.* Resolution of Disease Phenotypes Resulting from Multilocus Genomic Variation. *N. Engl. J. Med.* **376**, 21–31

- (2017).
9. Turro, E. *et al.* Whole-genome sequencing of patients with rare diseases in a national health system. *Nature* **583**, 96–102 (2020).
 10. Amberger, J. S., Bocchini, C. A., Scott, A. F. & Hamosh, A. OMIM.org: leveraging knowledge across phenotype–gene relationships. *Nucleic Acids Res.* **47**, D1038–D1043 (2019).
 11. DiStefano, M. T. *et al.* The Gene Curation Coalition: A global effort to harmonize gene–disease evidence resources. *Genet. Med. Off. J. Am. Coll. Med. Genet.* **24**, 1732–1742 (2022).
 12. Mungall, C. J. *et al.* The Monarch Initiative: an integrative data and analytic platform connecting phenotypes to genotypes across species. *Nucleic Acids Res.* **45**, D712–D722 (2017).
 13. Bamshad, M. J., Nickerson, D. A. & Chong, J. X. Mendelian Gene Discovery: Fast and Furious with No End in Sight. *Am. J. Hum. Genet.* **105**, 448–455 (2019).
 14. Dawood, M. *et al.* Using multiplexed functional data to reduce variant classification inequities in underrepresented populations. *Genome Med.* **16**, 143 (2024).
 15. Surl, D. *et al.* Clinician-Driven Reanalysis of Exome Sequencing Data From Patients With Inherited Retinal Diseases. *JAMA Netw. Open* **7**, e2414198 (2024).
 16. Seaby, E. G. *et al.* A gene pathogenicity tool ‘GenePy’ identifies missed biallelic diagnoses in the 100,000 Genomes Project. *Genet. Med. Off. J. Am. Coll. Med. Genet.* **26**, 101073 (2024).
 17. Liu, P. *et al.* Reanalysis of Clinical Exome Sequencing Data. *N. Engl. J. Med.* **380**, 2478–2480 (2019).
 18. Berger, S. I. *et al.* Increased diagnostic yield from negative whole genome–slice panels using automated reanalysis. *Clin. Genet.* **104**, 377–383 (2023).
 19. Wenger, A. M., Guturu, H., Bernstein, J. A. & Bejerano, G. Systematic reanalysis of clinical exome data yields additional diagnoses: implications for providers. *Genet. Med. Off. J. Am. Coll. Med. Genet.* **19**, 209–214 (2017).
 20. Zech, M. *et al.* Variants in Mitochondrial ATP Synthase Cause Variable Neurologic Phenotypes. *Ann. Neurol.* **91**, 225–237 (2022).
 21. Villamor-Payà, M. *et al.* De novo TLK1 and MDM1 mutations in a patient with a neurodevelopmental disorder and immunodeficiency. *iScience* **27**, 109984 (2024).
 22. Stegmann, J. D. *et al.* Bi-allelic variants in CELSR3 are implicated in central nervous system and urinary tract anomalies. *NPJ Genomic Med.* **9**, 18 (2024).
 23. Sobering, A. K. *et al.* Variants in PHF8 cause a spectrum of X-linked neurodevelopmental disorders and facial dysmorphism. *HGG Adv.* **3**, 100102 (2022).
 24. Serpieri, V. *et al.* SUFU haploinsufficiency causes a recognisable neurodevelopmental phenotype at the mild end of the Joubert syndrome spectrum. *J. Med. Genet.* **59**, 888–894 (2022).
 25. Sczakiel, H. L. *et al.* Broadening the phenotypic and molecular spectrum of FINCA syndrome: Biallelic NHLRC2 variants in 15 novel individuals. *Eur. J. Hum. Genet. EJHG* **31**, 905–917 (2023).
 26. Scala, M. *et al.* Variant-specific changes in RAC3 function disrupt corticogenesis in neurodevelopmental phenotypes. *Brain J. Neurol.* **145**, 3308–3327 (2022).
 27. Salinas, S. A. *et al.* An ELF4 hypomorphic variant results in NK cell deficiency. *JCI Insight* **7**, e155481 (2022).
 28. Saffari, A. *et al.* The clinical and genetic spectrum of autosomal-recessive TOR1A-related disorders. *Brain J. Neurol.* **146**, 3273–3288 (2023).
 29. Ren, X. *et al.* Increased TBX6 gene dosages induce congenital cervical vertebral malformations in humans and mice. *J. Med. Genet.* **57**, 371–379 (2020).
 30. Qian, X. *et al.* Loss of non-motor kinesin KIF26A causes congenital brain malformations via dysregulated neuronal migration and axonal growth as well as apoptosis. *Dev. Cell* **57**, 2381–2396.e13 (2022).
 31. Pérez Baca, M. D. R. *et al.* Haploinsufficiency of ZFH3, encoding a key player in neuronal development, causes syndromic intellectual disability. *Am. J. Hum. Genet.* **111**, 509–528 (2024).
 32. Penon-Portmann, M. *et al.* De novo heterozygous variants in SLC30A7 are a candidate cause for Joubert syndrome. *Am. J. Med. Genet. A.* **188**, 2360–2366 (2022).
 33. Muntadas, J. A. *et al.* Congenital myasthenic syndrome secondary to pathogenic variants in the SLC5A7 gene: report of two cases. *BMC Med. Genomics* **17**, 207 (2024).
 34. Münch, J. *et al.* Biallelic pathogenic variants in roundabout guidance receptor 1 associate with syndromic congenital anomalies of the kidney and urinary tract. *Kidney Int.* **101**, 1039–1053 (2022).
 35. Mubungu, G. *et al.* Clinical presentation and evolution of Xia-Gibbs syndrome due to p.Gly375ArgfsTer3 variant in a patient from DR Congo (Central Africa). *Am. J. Med. Genet. A.* **185**, 990–994 (2021).
 36. Morales-Rosado, J. A. *et al.* Bi-allelic variants in HMGCR cause an autosomal-recessive progressive limb-girdle muscular dystrophy. *Am. J. Hum. Genet.* **110**, 989–997 (2023).
 37. Mo, A. *et al.* Early-Onset and Severe Complex Hereditary Spastic Paraplegia Caused by De Novo Variants in SPAST. *Mov. Disord. Off. J. Mov. Disord. Soc.* **37**, 2440–2446 (2022).
 38. Mitani, T. *et al.* High prevalence of multilocus pathogenic variation in neurodevelopmental disorders in the Turkish population. *Am. J. Hum. Genet.* **108**, 1981–2005 (2021).
 39. Miller, D. E. *et al.* Targeted long-read sequencing identifies missing pathogenic variants in unsolved Werner syndrome cases. *J. Med. Genet.* **59**, 1087–1094 (2022).
 40. Miller, D. E. *et al.* Targeted Long-Read Sequencing Identifies a Retrotransposon Insertion as a Cause of Altered GNAS Exon A/B Methylation in a Family With Autosomal Dominant Pseudohypoparathyroidism Type 1b (PHP1B). *J. Bone Miner. Res. Off. J. Am. Soc. Bone Miner. Res.* **37**, 1711–1719 (2022).
 41. Maroofian, R. *et al.* Biallelic variants in SLC4A10 encoding a sodium-dependent bicarbonate transporter lead to a neurodevelopmental disorder. *Genet. Med. Off. J. Am. Coll. Med. Genet.* **26**, 101034 (2024).
 42. María Del Rocío, P. B. *et al.* Loss-of-function of the Zinc Finger Homeobox 4 (ZFX4) gene underlies a neurodevelopmental

- disorder. *MedRxiv Prepr. Serv. Health Sci.* 2024.08.07.24311381 (2024) doi:10.1101/2024.08.07.24311381.
43. Marafi, D. *et al.* A reverse genetics and genomics approach to gene paralog function and disease: Myokymia and the juxtapanaradome. *Am. J. Hum. Genet.* **109**, 1713–1723 (2022).
44. Mao, K. *et al.* FOXI3 pathogenic variants cause one form of craniofacial microsomia. *Nat. Commun.* **14**, 2026 (2023).
45. Liu, Z. *et al.* Hemizygous variants in protein phosphatase 1 regulatory subunit 3F (PPP1R3F) are associated with a neurodevelopmental disorder characterized by developmental delay, intellectual disability and autistic features. *Hum. Mol. Genet.* **32**, 2981–2995 (2023).
46. Liu, Q. *et al.* The Genetic Landscape of Familial Pulmonary Fibrosis. *Am. J. Respir. Crit. Care Med.* **207**, 1345–1357 (2023).
47. Lecca, M. *et al.* Bi-allelic variants in the ESAM tight-junction gene cause a neurodevelopmental disorder associated with fetal intracranial hemorrhage. *Am. J. Hum. Genet.* **110**, 681–690 (2023).
48. Laurent, S. *et al.* Molecular characterization of pathogenic OTOA gene conversions in hearing loss patients. *Hum. Mutat.* **42**, 373–377 (2021).
49. Khayat, M. M. *et al.* AHDC1 missense mutations in Xia-Gibbs syndrome. *HGG Adv.* **2**, 100049 (2021).
50. Keehan, L. *et al.* Wide range of phenotypic severity in individuals with late truncations unique to the predominant CDKL5 transcript in the brain. *Am. J. Med. Genet. A.* **188**, 3516–3524 (2022).
51. Kaiyrzhanov, R. *et al.* Bi-allelic ACBD6 variants lead to a neurodevelopmental syndrome with progressive and complex movement disorders. *Brain J. Neurol.* **147**, 1436–1456 (2024).
52. Kaiyrzhanov, R. *et al.* GGPS1-associated muscular dystrophy with and without hearing loss. *Ann. Clin. Transl. Neurol.* **9**, 1465–1474 (2022).
53. Huang, Y. *et al.* The recurrent de novo c.2011C>T missense variant in MTSS2 causes syndromic intellectual disability. *Am. J. Hum. Genet.* **109**, 1923–1931 (2022).
54. Hisama, F. M. *et al.* Caspase 5 depletion is linked to hyper-inflammatory response and progeroid syndrome. *GeroScience* **46**, 2771–2775 (2024).
55. Hijazi, H. *et al.* TCEAL1 loss-of-function results in an X-linked dominant neurodevelopmental syndrome and drives the neurological disease trait in Xq22.2 deletions. *Am. J. Hum. Genet.* **109**, 2270–2282 (2022).
56. Herman, I. *et al.* Quantitative dissection of multilocus pathogenic variation in an Egyptian infant with severe neurodevelopmental disorder resulting from multiple molecular diagnoses. *Am. J. Med. Genet. A.* **188**, 735–750 (2022).
57. Grammatikopoulos, T. *et al.* Liver Disease and Risk of Hepatocellular Carcinoma in Children With Mutations in TALDO1. *Hepatol. Commun.* **6**, 473–479 (2022).
58. Gourgas, O. *et al.* Specific heterozygous variants in MGP lead to endoplasmic reticulum stress and cause spondyloepiphyseal dysplasia. *Nat. Commun.* **14**, 7054 (2023).
59. Gofin, Y. *et al.* Delineation of a novel neurodevelopmental syndrome associated with PAX5 haploinsufficiency. *Hum. Mutat.* **43**, 461–470 (2022).
60. Furia, F. *et al.* The phenotypic and genotypic spectrum of individuals with mono- or biallelic ANK3 variants. *Clin. Genet.* **106**, 574–584 (2024).
61. Frost, F. G. *et al.* Bi-allelic SNAPC4 variants dysregulate global alternative splicing and lead to neuroregression and progressive spastic paraparesis. *Am. J. Hum. Genet.* **110**, 663–680 (2023).
62. Foley, A. R. *et al.* The recurrent deep intronic pseudoexon-inducing variant COL6A1 c.930+189C>T results in a consistently severe phenotype of COL6-related dystrophy: Towards clinical trial readiness for splice-modulating therapy. *MedRxiv Prepr. Serv. Health Sci.* 2024.03.29.24304673 (2024) doi:10.1101/2024.03.29.24304673.
63. Feng, X. *et al.* Core planar cell polarity genes VANGL1 and VANGL2 in predisposition to congenital vertebral malformations. *Proc. Natl. Acad. Sci. U. S. A.* **121**, e2310283121 (2024).
64. Fasham, J. *et al.* Elucidating the clinical spectrum and molecular basis of HYAL2 deficiency. *Genet. Med. Off. J. Am. Coll. Med. Genet.* **24**, 631–644 (2022).
65. Faqeih, E. A. *et al.* Biallelic variants in HECT E3 paralogs, HECTD4 and UBE3C, encoding ubiquitin ligases cause neurodevelopmental disorders that overlap with Angelman syndrome. *Genet. Med. Off. J. Am. Coll. Med. Genet.* **25**, 100323 (2023).
66. Elbendary, H. M. *et al.* Novel LSS variants in alopecia and intellectual disability syndrome: New case report and clinical spectrum of LSS-related rare disease traits. *Clin. Genet.* **104**, 344–349 (2023).
67. Duan, R. *et al.* Biallelic missense variants in COG3 cause a congenital disorder of glycosylation with impairment of retrograde vesicular trafficking. *J. Inherit. Metab. Dis.* **46**, 1195–1205 (2023).
68. Duan, R. *et al.* Developmental genomics of limb malformations: Allelic series in association with gene dosage effects contribute to the clinical variability. *HGG Adv.* **3**, 100132 (2022).
69. Dias, K.-R. *et al.* Narrowing the diagnostic gap: Genomes, epigenomes, long-read sequencing, and health economic analyses in an exome-negative intellectual disability cohort. *Genet. Med. Off. J. Am. Coll. Med. Genet.* **26**, 101076 (2024).
70. Dawood, M. *et al.* A biallelic frameshift indel in PPP1R35 as a cause of primary microcephaly. *Am. J. Med. Genet. A.* **191**, 794–804 (2023).
71. Dardas, Z. *et al.* NODAL variants are associated with a continuum of laterality defects from simple D-transposition of the great arteries to heterotaxy. *Genome Med.* **16**, 53 (2024).
72. Copeland, I. *et al.* Exome sequencing implicates ancestry-related Mendelian variation at SYNE1 in childhood-onset essential hypertension. *JCI Insight* **9**, e172152 (2024).
73. Cingöz, S. *et al.* Novel biallelic variants affecting the OTU domain of the gene OTUD6B associate with severe intellectual disability syndrome and molecular dynamics simulations. *Eur. J. Med. Genet.* **65**, 104497 (2022).
74. Chong, J. X. *et al.* Variants in ACTC1 underlie distal arthrogyposis accompanied by congenital heart defects. *HGG Adv.* **4**, 100213 (2023).
75. Chen, Y. *et al.* De novo variants in the RNU4-2 snRNA cause a frequent neurodevelopmental syndrome. *Nature* **632**, 832–840

- (2024).
76. Cetica, V. *et al.* Clinical and molecular characterization of patients with YWHAG-related epilepsy. *Epilepsia* **65**, 1439–1450 (2024).
 77. Caron, V. *et al.* Clinical and functional heterogeneity associated with the disruption of retinoic acid receptor beta. *Genet. Med. Off. J. Am. Coll. Med. Genet.* **25**, 100856 (2023).
 78. Calame, D. G. *et al.* Cation leak through the ATP1A3 pump causes spasticity and intellectual disability. *Brain J. Neurol.* **146**, 3162–3171 (2023).
 79. Calame, D. G. *et al.* Biallelic Variants in the Ectonucleotidase ENTPD1 Cause a Complex Neurodevelopmental Disorder with Intellectual Disability, Distinct White Matter Abnormalities, and Spastic Paraplegia. *Ann. Neurol.* **92**, 304–321 (2022).
 80. Calame, D. G. *et al.* Monoallelic variation in DHX9, the gene encoding the DExH-box helicase DHX9, underlies neurodevelopment disorders and Charcot-Marie-Tooth disease. *Am. J. Hum. Genet.* **110**, 1394–1413 (2023).
 81. Brown, G. J. *et al.* TLR7 gain-of-function genetic variation causes human lupus. *Nature* **605**, 349–356 (2022).
 82. Brooks, D. *et al.* Heterozygous MAP3K20 variants cause ectodermal dysplasia, craniosynostosis, sensorineural hearing loss, and limb anomalies. *Hum. Genet.* **143**, 279–291 (2024).
 83. Boschann, F. *et al.* Biallelic variants in ADAMTS15 cause a novel form of distal arthrogyposis. *Genet. Med. Off. J. Am. Coll. Med. Genet.* **24**, 2187–2193 (2022).
 84. Borroto, M. C. *et al.* A Genotype/Phenotype Study of KDM5B-Associated Disorders Suggests a Pathogenic Effect of Dominantly Inherited Missense Variants. *Genes* **15**, 1033 (2024).
 85. Blue, E. E. *et al.* Exome sequencing identifies novel genes underlying primary congenital glaucoma in the National Birth Defects Prevention Study. *Birth Defects Res.* **116**, e2384 (2024).
 86. Berger, S. I., Miller, I. & Tochen, L. Recessive GCH1 Deficiency Causing DOPA-Responsive Dystonia Diagnosed by Reported Negative Exome. *Pediatrics* **149**, e2021052886 (2022).
 87. Bassani, S. *et al.* Variant-specific pathophysiological mechanisms of AFF3 differently influence transcriptome profiles. *Genome Med.* **16**, 72 (2024).
 88. Banks, E. *et al.* Loss of symmetric cell division of apical neural progenitors drives DENND5A-related developmental and epileptic encephalopathy. *Nat. Commun.* **15**, 7239 (2024).
 89. Ansari, M. *et al.* Heterozygous loss-of-function SMC3 variants are associated with variable growth and developmental features. *HGG Adv.* **5**, 100273 (2024).
 90. Almannai, M. *et al.* Expanding the phenotype of PPP1R21-related neurodevelopmental disorder. *Clin. Genet.* **105**, 620–629 (2024).
 91. Almannai, M. *et al.* El-Hattab-Alkuraya syndrome caused by biallelic WDR45B pathogenic variants: Further delineation of the phenotype and genotype. *Clin. Genet.* **101**, 530–540 (2022).
 92. Al-Kouatly, H. B. *et al.* High diagnosis rate for nonimmune hydrops fetalis with prenatal clinical exome from the Hydrops-Yielding Diagnostic Results of Prenatal Sequencing (HYDROPS) Study. *Genet. Med. Off. J. Am. Coll. Med. Genet.* **23**, 1325–1333 (2021).
 93. Ganesh, V. S. *et al.* Neurodevelopmental Disorder Caused by Deletion of CHASERR, a lncRNA Gene. *N. Engl. J. Med.* **391**, 1511–1518 (2024).
 94. Rai, A. *et al.* Genomic rare variant mechanisms for congenital cardiac laterality defect: A digenic model approach. 2024.11.19.24317385 Preprint at <https://doi.org/10.1101/2024.11.19.24317385> (2024).
 95. Grasberger, H. *et al.* STR mutations on chromosome 15q cause thyrotropin resistance by activating a primate-specific enhancer of MIR7-2/MIR1179. *Nat. Genet.* **56**, 877–888 (2024).
 96. Jolly, A. *et al.* Rare variant enrichment analysis supports GREB1L as a contributory driver gene in the etiology of Mayer-Rokitansky-Küster-Hauser syndrome. *HGG Adv.* **4**, 100188 (2023).
 97. Mori, T. *et al.* CFP47 is Implicated in X-Linked Polycystic Kidney Disease. *Kidney Int. Rep.* **9**, 3580–3591 (2024).
 98. Ma, M. *et al.* De novo variants in PLAG1 are associated with hearing impairment, ocular pathology, and cardiac defects. *eLife* **13**, (2024).
 99. Calame, D. G. *et al.* Biallelic variation in the choline and ethanolamine transporter FLVCR1 underlies a severe developmental disorder spectrum. *Genet. Med. Off. J. Am. Coll. Med. Genet.* 101273 (2024) doi:10.1016/j.gim.2024.101273.
 100. Weisburd, B. *et al.* Diagnosing missed cases of spinal muscular atrophy in genome, exome, and panel sequencing datasets. *MedRxiv Prepr. Serv. Health Sci.* 2024.02.11.24302646 (2024) doi:10.1101/2024.02.11.24302646.
 101. Guo, M. H. *et al.* Inferring compound heterozygosity from large-scale exome sequencing data. *Nat. Genet.* **56**, 152–161 (2024).
 102. Gudmundsson, S. *et al.* Variant interpretation using population databases: Lessons from gnomAD. *Hum. Mutat.* **43**, 1012–1030 (2022).
 103. Lemire, G. *et al.* Exome copy number variant detection, analysis, and classification in a large cohort of families with undiagnosed rare genetic disease. *Am. J. Hum. Genet.* **111**, 863–876 (2024).
 104. Du, H. *et al.* HMZDupFinder: a robust computational approach for detecting intragenic homozygous duplications from exome sequencing data. *Nucleic Acids Res.* **52**, e18 (2024).
 105. Babadi, M. *et al.* GATK-gCNV enables the discovery of rare copy number variants from exome sequencing data. *Nat. Genet.* **55**, 1589–1597 (2023).
 106. Du, H. *et al.* VizCNV: An integrated platform for concurrent phased BAF and CNV analysis with trio genome sequencing data. Preprint at <https://doi.org/10.1101/2024.10.27.620363> (2024).
 107. Wojcik, M. H. *et al.* Beyond the exome: What's next in diagnostic testing for Mendelian conditions. *Am. J. Hum. Genet.* **110**, 1229–1248 (2023).
 108. All of Us Research Program Genomics Investigators. Genomic data in the All of Us Research Program. *Nature* **627**, 340–346 (2024).
 109. Li, S., Carss, K. J., Halldorsson, B. V., Cortes, A., & UK Biobank Whole-Genome Sequencing Consortium. Whole-genome sequencing of half-a-million UK Biobank participants. Preprint at <https://doi.org/10.1101/2023.12.06.23299426> (2023).

110. Wenger, T. L. *et al.* SeqFirst: Building equity access to a precise genetic diagnosis in critically ill newborns. Preprint at <https://doi.org/10.1101/2024.09.30.24314516> (2024).
111. Wojcik, M. H. *et al.* Genome Sequencing for Diagnosing Rare Diseases. *N. Engl. J. Med.* **390**, 1985–1997 (2024).
112. Saad, A. K. *et al.* Biallelic in-frame deletion in TRAPPC4 in a family with developmental delay and cerebellar atrophy. *Brain J. Neurol.* **143**, e83 (2020).
113. Bruels, C. C. *et al.* Diagnostic capabilities of nanopore long-read sequencing in muscular dystrophy. *Ann. Clin. Transl. Neurol.* **9**, 1302–1309 (2022).
114. Chen, X. *et al.* Genome-wide profiling of highly similar paralogous genes using HiFi sequencing. 2024.04.19.590294 Preprint at <https://doi.org/10.1101/2024.04.19.590294> (2024).
115. Negi, S. *et al.* Advancing long-read nanopore genome assembly and accurate variant calling for rare disease detection. *MedRxiv Prepr. Serv. Health Sci.* 2024.08.22.24312327 (2024) doi:10.1101/2024.08.22.24312327.
116. Mahmoud, M. *et al.* Closing the gap: Solving complex medically relevant genes at scale. *MedRxiv Prepr. Serv. Health Sci.* 2024.03.14.24304179 (2024) doi:10.1101/2024.03.14.24304179.
117. Kronenberg, Z. *et al.* The Platinum Pedigree: A long-read benchmark for genetic variants. Preprint at <https://doi.org/10.1101/2024.10.02.616333> (2024).
118. Majidian, S., Agostinho, D. P., Chin, C.-S., Sedlazeck, F. J. & Mahmoud, M. Genomic variant benchmark: if you cannot measure it, you cannot improve it. *Genome Biol.* **24**, 221 (2023).
119. LoTempio, J., Delot, E. & Vilain, E. Benchmarking long-read genome sequence alignment tools for human genomics applications. *PeerJ* **11**, e16515 (2023).
120. Liu, Z. *et al.* Towards accurate and reliable resolution of structural variants for clinical diagnosis. *Genome Biol.* **23**, 68 (2022).
121. Gustafson, J. A. *et al.* High-coverage nanopore sequencing of samples from the 1000 Genomes Project to build a comprehensive catalog of human genetic variation. *Genome Res.* **34**, 2061–2073 (2024).
122. Zheng, X. *et al.* STIX: Long-reads based Accurate Structural Variation Annotation at Population Scale. Preprint at <https://doi.org/10.1101/2024.09.30.615931> (2024).
123. Smolka, M. *et al.* Detection of mosaic and population-level structural variants with Sniffles2. *Nat. Biotechnol.* **42**, 1571–1580 (2024).
124. Sedlazeck, F. J. *et al.* Accurate detection of complex structural variations using single-molecule sequencing. *Nat. Methods* **15**, 461–468 (2018).
125. Bilgrav Saether, K. *et al.* Leveraging the T2T assembly to resolve rare and pathogenic inversions in reference genome gaps. *Genome Res.* (2024) doi:10.1101/gr.279346.124.
126. Fu, Y. *et al.* MethPhaser: methylation-based long-read haplotype phasing of human genomes. *Nat. Commun.* **15**, 5327 (2024).
127. LaFlamme, C. W. *et al.* Diagnostic utility of DNA methylation analysis in genetically unsolved pediatric epilepsies and CHD2 epismutation refinement. *Nat. Commun.* **15**, 6524 (2024).
128. Stergachis, A. B., Debo, B. M., Haugen, E., Churchman, L. S. & Stamatoyannopoulos, J. A. Single-molecule regulatory architectures captured by chromatin fiber sequencing. *Science* **368**, 1449–1454 (2020).
129. Jha, A. *et al.* DNA-m6A calling and integrated long-read epigenetic and genetic analysis with fibertools. *Genome Res.* gr.279095.124 (2024) doi:10.1101/gr.279095.124.
130. Vollger, M. R. *et al.* Synchronized long-read genome, methylome, epigenome, and transcriptome for resolving a Mendelian condition. *BioRxiv Prepr. Serv. Biol.* 2023.09.26.559521 (2023) doi:10.1101/2023.09.26.559521.
131. Carvalho, C. M. B. *et al.* Interchromosomal template-switching as a novel molecular mechanism for imprinting perturbations associated with Temple syndrome. *Genome Med.* **11**, 25 (2019).
132. Smail, C. *et al.* Integration of rare expression outlier-associated variants improves polygenic risk prediction. *Am. J. Hum. Genet.* **109**, 1055–1064 (2022).
133. Li, T. *et al.* The functional impact of rare variation across the regulatory cascade. *Cell Genomics* **3**, 100401 (2023).
134. Brechtmann, F. *et al.* OUTRIDER: A Statistical Method for Detecting Aberrantly Expressed Genes in RNA Sequencing Data. *Am. J. Hum. Genet.* **103**, 907–917 (2018).
135. Montgomery, S. B., Bernstein, J. A. & Wheeler, M. T. Toward transcriptomics as a primary tool for rare disease investigation. *Cold Spring Harb. Mol. Case Stud.* **8**, a006198 (2022).
136. Jensen, T. D. *et al.* Integration of transcriptomics and long-read genomics prioritizes structural variants in rare disease. *MedRxiv Prepr. Serv. Health Sci.* 2024.03.22.24304565 (2024) doi:10.1101/2024.03.22.24304565.
137. Pais, L. S. *et al.* seqr: A web-based analysis and collaboration tool for rare disease genomics. *Hum. Mutat.* **43**, 698–707 (2022).
138. Yamamoto, S. *et al.* A drosophila genetic resource of mutants to study mechanisms underlying human genetic diseases. *Cell* **159**, 200–214 (2014).
139. Marderstein, A. R. *et al.* Single-cell multi-omics map of human fetal blood in Down syndrome. *Nature* **634**, 104–112 (2024).
140. Ober-Reynolds, B. *et al.* Integrated single-cell chromatin and transcriptomic analyses of human scalp identify gene-regulatory programs and critical cell types for hair and skin diseases. *Nat. Genet.* **55**, 1288–1300 (2023).
141. Corces, M. R. *et al.* Single-cell epigenomic analyses implicate candidate causal variants at inherited risk loci for Alzheimer's and Parkinson's diseases. *Nat. Genet.* **52**, 1158–1168 (2020).
142. Chen, A. F. *et al.* NEAT-seq: simultaneous profiling of intra-nuclear proteins, chromatin accessibility and gene expression in single cells. *Nat. Methods* **19**, 547–553 (2022).
143. Kalef-Ezra, E. *et al.* Single-cell somatic copy number variants in brain using different amplification methods and reference genomes. *Commun. Biol.* **7**, 1288 (2024).
144. IGVF Consortium. Deciphering the impact of genomic variation on function. *Nature* **633**, 47–57 (2024).
145. Fowler, D. M. & Fields, S. Deep mutational scanning: a new style of protein science. *Nat. Methods* **11**, 801–807 (2014).

146. Melnikov, A. *et al.* Systematic dissection and optimization of inducible enhancers in human cells using a massively parallel reporter assay. *Nat. Biotechnol.* **30**, 271–277 (2012).
147. Scott, H. A. *et al.* A high throughput splicing assay to investigate the effect of variants of unknown significance on exon inclusion. Preprint at <https://doi.org/10.1101/2022.11.30.22282952> (2022).
148. Lacoste, J. *et al.* Pervasive mislocalization of pathogenic coding variants underlying human disorders. *Cell* S0092-8674(24)01021–3 (2024) doi:10.1016/j.cell.2024.09.003.
149. Sahni, N. *et al.* Widespread macromolecular interaction perturbations in human genetic disorders. *Cell* **161**, 647–660 (2015).
150. Yi, S. *et al.* Functional variomics and network perturbation: connecting genotype to phenotype in cancer. *Nat. Rev. Genet.* **18**, 395–410 (2017).
151. Tai, D. J. C. *et al.* Tissue- and cell-type-specific molecular and functional signatures of 16p11.2 reciprocal genomic disorder across mouse brain and human neuronal models. *Am. J. Hum. Genet.* **109**, 1789–1813 (2022).
152. Mohajeri, K. *et al.* Transcriptional and functional consequences of alterations to MEF2C and its topological organization in neuronal models. *Am. J. Hum. Genet.* **109**, 2049–2067 (2022).
153. Pampari, A. *et al.* Bias factorized, base-resolution deep learning models of chromatin accessibility reveal cis-regulatory sequence syntax, transcription factor footprints and regulatory variants. Zenodo <https://doi.org/10.5281/zenodo.10396047> (2023).
154. Lansdon, L. A. *et al.* Factors Affecting Migration to GRCh38 in Laboratories Performing Clinical Next-Generation Sequencing. *J. Mol. Diagn. JMD* **23**, 651–657 (2021).
155. Li, H. *et al.* Exome variant discrepancies due to reference-genome differences. *Am. J. Hum. Genet.* **108**, 1239–1250 (2021).
156. Ungar, R. A. *et al.* Impact of genome build on RNA-seq interpretation and diagnostics. *Am. J. Hum. Genet.* **111**, 1282–1300 (2024).
157. Behera, S. *et al.* FixItFelix: improving genomic analysis by fixing reference errors. *Genome Biol.* **24**, 31 (2023).
158. Wagner, J. *et al.* Curated variation benchmarks for challenging medically relevant autosomal genes. *Nat. Biotechnol.* **40**, 672–680 (2022).
159. Liao, W.-W. *et al.* A draft human pangenome reference. *Nature* **617**, 312–324 (2023).
160. Behera, S. *et al.* Comprehensive genome analysis and variant detection at scale using DRAGEN. *Nat. Biotechnol.* (2024) doi:10.1038/s41587-024-02382-1.
161. Chin, C.-S. *et al.* Multiscale analysis of pangenomes enables improved representation of genomic diversity for repetitive and clinically relevant genes. *Nat. Methods* **20**, 1213–1221 (2023).
162. Chin, C.-S. *et al.* A pan-genome approach to decipher variants in the highly complex tandem repeat of LPA. 2022.06.08.495395 Preprint at <https://doi.org/10.1101/2022.06.08.495395> (2022).
163. Köhler, S. *et al.* The Human Phenotype Ontology in 2021. *Nucleic Acids Res.* **49**, D1207–D1217 (2021).
164. Zhang, C. *et al.* Novel pathogenic variants and quantitative phenotypic analyses of Robinow syndrome: WNT signaling perturbation and phenotypic variability. *HGG Adv.* **3**, 100074 (2022).
165. Garcia, B. T. *et al.* Improving Automated Deep Phenotyping Through Large Language Models Using Retrieval Augmented Generation. Preprint at <https://doi.org/10.1101/2024.12.01.24318253> (2024).
166. Samocha, K. E. *et al.* A framework for the interpretation of de novo mutation in human disease. *Nat. Genet.* **46**, 944–950 (2014).
167. Lupski, J. R., Belmont, J. W., Boerwinkle, E. & Gibbs, R. A. Clan genomics and the complex architecture of human disease. *Cell* **147**, 32–43 (2011).
168. Maquat, L. E. Nonsense-mediated mRNA decay in mammals. *J. Cell Sci.* **118**, 1773–1776 (2005).
169. Maquat, L. E. Nonsense-mediated mRNA decay. *Curr. Biol. CB* **12**, R196–197 (2002).
170. Teran, N. A. *et al.* Nonsense-mediated decay is highly stable across individuals and tissues. *Am. J. Hum. Genet.* **108**, 1401–1408 (2021).
171. Coban-Akdemir, Z. *et al.* Identifying Genes Whose Mutant Transcripts Cause Dominant Disease Traits by Potential Gain-of-Function Alleles. *Am. J. Hum. Genet.* **103**, 171–187 (2018).
172. Valencia, A. M. *et al.* Landscape of mSWI/SNF chromatin remodeling complex perturbations in neurodevelopmental disorders. *Nat. Genet.* **55**, 1400–1412 (2023).
173. Paine, I. *et al.* Paralog Studies Augment Gene Discovery: DDX and DHX Genes. *Am. J. Hum. Genet.* **105**, 302–316 (2019).
174. Gillentine, M. A. *et al.* Rare deleterious mutations of HNRNP genes result in shared neurodevelopmental disorders. *Genome Med.* **13**, 63 (2021).
175. Ochoa, S. *et al.* A deep intronic splice-altering AIRE variant causes APECED syndrome through antisense oligonucleotide-targetable pseudoexon inclusion. *Sci. Transl. Med.* **16**, eadk0845 (2024).
176. Jaganathan, K. *et al.* Predicting Splicing from Primary Sequence with Deep Learning. *Cell* **176**, 535–548.e24 (2019).
177. Wu, N. *et al.* TBX6 null variants and a common hypomorphic allele in congenital scoliosis. *N. Engl. J. Med.* **372**, 341–350 (2015).
178. Lord, J. *et al.* Non-coding variants are a rare cause of recessive developmental disorders in trans with coding variants. *Genet. Med. Off. J. Am. Coll. Med. Genet.* 101249 (2024) doi:10.1016/j.gim.2024.101249.
179. Greene, D. *et al.* Mutations in the U4 snRNA gene RNU4-2 cause one of the most prevalent monogenic neurodevelopmental disorders. *Nat. Med.* **30**, 2165–2169 (2024).
180. Greene, D. *et al.* Mutations in the U2 snRNA gene RNU2-2P cause a severe neurodevelopmental disorder with prominent epilepsy. *MedRxiv Prepr. Serv. Health Sci.* 2024.09.03.24312863 (2024) doi:10.1101/2024.09.03.24312863.
181. Nava, C. *et al.* Dominant variants in major spliceosome U4 and U5 small nuclear RNA genes cause neurodevelopmental disorders through splicing disruption. 2024.10.07.24314689 Preprint at <https://doi.org/10.1101/2024.10.07.24314689> (2024).
182. Töpf, A. *et al.* Digenic inheritance involving a muscle-specific protein kinase and the giant titin protein causes a skeletal muscle myopathy. *Nat. Genet.* **56**, 395–407 (2024).

183. Gifford, C. A. *et al.* Oligogenic inheritance of a human heart disease involving a genetic modifier. *Science* **364**, 865–870 (2019).
184. Bozkurt-Yozgatli, T. *et al.* Multilocus pathogenic variants contribute to intrafamilial clinical heterogeneity: a retrospective study of sibling pairs with neurodevelopmental disorders. *BMC Med. Genomics* **17**, 85 (2024).
185. Abell, N. S. *et al.* Multiple causal variants underlie genetic associations in humans. *Science* **375**, 1247–1254 (2022).
186. Lupski, J. R. Biology in balance: human diploid genome integrity, gene dosage, and genomic medicine. *Trends Genet. TIG* **38**, 554–571 (2022).
187. Liu, P. *et al.* An Organismal CNV Mutator Phenotype Restricted to Early Human Development. *Cell* **168**, 830–842.e7 (2017).
188. Du, H. *et al.* The multiple de novo copy number variant (MdnCNV) phenomenon presents with peri-zygotic DNA mutational signatures and multilocus pathogenic variation. *Genome Med.* **14**, 122 (2022).
189. Logsdon, G. A. *et al.* Complex genetic variation in nearly complete human genomes. *BioRxiv Prepr. Serv. Biol.* 2024.09.24.614721 (2024) doi:10.1101/2024.09.24.614721.
190. Ebert, P. *et al.* Haplotype-resolved diverse human genomes and integrated analysis of structural variation. *Science* **372**, eabf7117 (2021).
191. Collins, R. L. *et al.* A structural variation reference for medical and population genetics. *Nature* **581**, 444–451 (2020).
192. Byrska-Bishop, M. *et al.* High-coverage whole-genome sequencing of the expanded 1000 Genomes Project cohort including 602 trios. *Cell* **185**, 3426–3440.e19 (2022).
193. Grochowski, C. M. *et al.* Inverted triplications formed by iterative template switches generate structural variant diversity at genomic disorder loci. *Cell Genomics* **4**, 100590 (2024).
194. Dardas, Z. *et al.* Genomic Balancing Act: deciphering DNA rearrangements in the complex chromosomal aberration involving 5p15.2, 2q31.1, and 18q21.32. *Eur. J. Hum. Genet. EJHG* (2024) doi:10.1038/s41431-024-01680-1.
195. Bilgrav Saether, K. *et al.* Mind the gap: the relevance of the genome reference to resolve rare and pathogenic inversions. *MedRxiv Prepr. Serv. Health Sci.* 2024.04.22.24305780 (2024) doi:10.1101/2024.04.22.24305780.
196. Jakubosky, D. *et al.* Discovery and quality analysis of a comprehensive set of structural variants and short tandem repeats. *Nat. Commun.* **11**, 2928 (2020).
197. Jakubosky, D. *et al.* Properties of structural variants and short tandem repeats associated with gene expression and complex traits. *Nat. Commun.* **11**, 2927 (2020).
198. Dolzhenko, E. *et al.* Characterization and visualization of tandem repeats at genome scale. *Nat. Biotechnol.* **42**, 1606–1614 (2024).
199. English, A. C. *et al.* Analysis and benchmarking of small and large genomic variants across tandem repeats. *Nat. Biotechnol.* (2024) doi:10.1038/s41587-024-02225-z.
200. Dolzhenko, E. *et al.* REViewer: haplotype-resolved visualization of read alignments in and around tandem repeats. *Genome Med.* **14**, 84 (2022).
201. Behera, S. *et al.* Identification of allele-specific KIV-2 repeats and impact on Lp(a) measurements for cardiovascular disease risk. *BMC Med. Genomics* **17**, 255 (2024).
202. Weisburd, B., Tiao, G. & Rehm, H. L. Insights from a genome-wide truth set of tandem repeat variation. *BioRxiv Prepr. Serv. Biol.* 2023.05.05.539588 (2023) doi:10.1101/2023.05.05.539588.
203. Cui, Y. *et al.* A genome-wide spectrum of tandem repeat expansions in 338,963 humans. *Cell* **187**, 2336–2341.e5 (2024).
204. Weisburd, B. *et al.* Defining a tandem repeat catalog and variation clusters for genome-wide analyses and population databases. Preprint at <https://doi.org/10.1101/2024.10.04.615514> (2024).
205. Wang, Q. *et al.* Landscape of multi-nucleotide variants in 125,748 human exomes and 15,708 genomes. *Nat. Commun.* **11**, 2539 (2020).
206. Srinivasan, S. *et al.* Misannotated Multi-Nucleotide Variants in Public Cancer Genomics Datasets Lead to Inaccurate Mutation Calls with Significant Implications. *Cancer Res.* **81**, 282–288 (2021).
207. Campbell, I. M. *et al.* Multiallelic Positions in the Human Genome: Challenges for Genetic Analyses. *Hum. Mutat.* **37**, 231–234 (2016).
208. Singer-Berk, M. *et al.* Advanced variant classification framework reduces the false positive rate of predicted loss-of-function variants in population sequencing data. *Am. J. Hum. Genet.* **110**, 1496–1508 (2023).
209. Lindeboom, R. G. H., Supek, F. & Lehner, B. The rules and impact of nonsense-mediated mRNA decay in human cancers. *Nat. Genet.* **48**, 1112–1118 (2016).
210. Lindeboom, R. G. H., Vermeulen, M., Lehner, B. & Supek, F. The impact of nonsense-mediated mRNA decay on genetic disease, gene editing and cancer immunotherapy. *Nat. Genet.* **51**, 1645–1651 (2019).
211. Torene, R. I. *et al.* Systematic analysis of variants escaping nonsense-mediated decay uncovers candidate Mendelian diseases. *Am. J. Hum. Genet.* **111**, 70–81 (2024).
212. Potter, A. S. *et al.* Rare Variant in MRC2 Associated With Familial Supraventricular Tachycardia and Wolff-Parkinson-White Syndrome. *Circ. Genomic Precis. Med.* **17**, e004614 (2024).
213. Gudmundsson, S. *et al.* Exploring penetrance of clinically relevant variants in over 800,000 humans from the Genome Aggregation Database. 2024.06.12.593113 Preprint at <https://doi.org/10.1101/2024.06.12.593113> (2024).
214. 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. Off. J. Am. Coll. Med. Genet.* **17**, 405–424 (2015).
215. Rehm, H. L. *et al.* The landscape of reported VUS in multi-gene panel and genomic testing: Time for a change. *Genet. Med. Off. J. Am. Coll. Med. Genet.* **25**, 100947 (2023).
216. Fayer, S. *et al.* Closing the gap: Systematic integration of multiplexed functional data resolves variants of uncertain significance in BRCA1, TP53, and PTEN. *Am. J. Hum. Genet.* **108**, 2248–2258 (2021).

217. Young, J. L. *et al.* Beyond race: Recruitment of diverse participants in clinical genomics research for rare disease. *Front. Genet.* **13**, 949422 (2022).
218. Wojcik, M. H. *et al.* Rare diseases, common barriers: disparities in pediatric clinical genetics outcomes. *Pediatr. Res.* **93**, 110–117 (2023).
219. Serrano, J. G. *et al.* Advancing Understanding of Inequities in Rare Disease Genomics. *Clin. Ther.* **45**, 745–753 (2023).
220. D'Angelo, C. S. *et al.* Barriers and Considerations for Diagnosing Rare Diseases in Indigenous Populations. *Front. Pediatr.* **8**, 579924 (2020).
221. Savage, S. K. *et al.* Using a chat-based informed consent tool in large-scale genomic research. *J. Am. Med. Inform. Assoc. JAMIA* **31**, 472–478 (2024).
222. Stark, Z. *et al.* A call to action to scale up research and clinical genomic data sharing. *Nat. Rev. Genet.* (2024). doi:10.1038/s41576-024-00776-0.
223. Rehm, H. L. Time to make rare disease diagnosis accessible to all. *Nat. Med.* **28**, 241–242 (2022).
224. Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and stewardship. *Sci. Data* **3**, 160018 (2016).
225. Seaby, E. G. *et al.* A Panel-Agnostic Strategy 'HiPPo' Improves Diagnostic Efficiency in the UK Genomic Medicine Service. *Healthc. Basel Switz.* **11**, 3179 (2023).
226. Chong, J. X. *et al.* Considerations for reporting variants in novel candidate genes identified during clinical genomic testing. *Genet. Med. Off. J. Am. Coll. Med. Genet.* **26**, 101199 (2024).
227. Philippakis, A. A. *et al.* The Matchmaker Exchange: a platform for rare disease gene discovery. *Hum. Mutat.* **36**, 915–921 (2015).
228. Boycott, K. M., Azzariti, D. R., Hamosh, A. & Rehm, H. L. Seven years since the launch of the Matchmaker Exchange: The evolution of genomic matchmaking. *Hum. Mutat.* **43**, 659–667 (2022).
229. MyGene2, NHGRI/NHLBI University of Washington-Center for Mendelian Genomics (UW-CMG), Seattle, WA. <https://mygene2.org>.
230. Landrum, M. J. *et al.* ClinVar: improving access to variant interpretations and supporting evidence. *Nucleic Acids Res.* **46**, D1062–D1067 (2018).
231. Landrum, M. J. *et al.* ClinVar: public archive of relationships among sequence variation and human phenotype. *Nucleic Acids Res.* **42**, D980–985 (2014).
232. Wohler, E. *et al.* PhenoDB, GeneMatcher and VariantMatcher, tools for analysis and sharing of sequence data. *Orphanet J. Rare Dis.* **16**, 365 (2021).
233. Rodrigues, E. da S. *et al.* Variant-level matching for diagnosis and discovery: Challenges and opportunities. *Hum. Mutat.* **43**, 782–790 (2022).
234. Wagner, A. H. *et al.* The GA4GH Variation Representation Specification: A computational framework for variation representation and federated identification. *Cell Genomics* **1**, 100027 (2021).
235. Pawliczek, P. *et al.* ClinGen Allele Registry links information about genetic variants. *Hum. Mutat.* **39**, 1690–1701 (2018).
236. GREGoR Variant Browser. <https://variants.gregorconsortium.org/>.

Acknowledgements:

We would like to acknowledge all patient participants and their families. We would also like to acknowledge the expansive set of collaborators including clinical providers, analysts and rare disease researchers. This work was supported by the NIH NHGRI GREGoR Consortium (U01HG011758, U01HG011755, U01HG011762, U01HG011745, U01HG011744, U24HG011746).

Author contributions:

Conceptualization, Methodology: MD, BH, MMW, RAU, JLT, LW, SB, JAB, JXC, ECD, EEE, RAG, JRL, AS, MET, AHW, CLW, CW, MTW, CMBC, CAG, SM, DEM, HLR, FJS, EV, AODL, JEP, LHC, MJB, SBM; Data curation: MD, BH, MMW, SBM; Investigation: SB, JAB, JXC, ECD, EEE, RAG, JRL, AS, MET, AHW, CLW, CW, MTW, CMBC, CAG, SM, DEM, HLR, FJS, EV, AODL, JEP, LHC, MJB, SBM; Data curation: MD, BH, MMW, SBM; Project administration: SBM; Writing - original draft: MD, BH, MMW, RAU, JLT, LW, CMBC, CAG, SM, DEM, HLR, FJS, EV, AODL, JEP, LHC, MJB, SBM; All authors contributed to developing and enacting the vision and goals of the GREGoR Consortium, reviewing, and editing the manuscript.

Competing Interests:

All authors of this manuscript are funded by the NIH and NHGRI. JRL has stock ownership in 23andMe, is a paid consultant for Regeneron Genetics Center, and is a co-inventor on multiple U.S. and European patents related to molecular diagnostics for inherited neuropathies, eye diseases, and bacterial genomic fingerprinting. JRL serves on the Scientific Advisory Board of Baylor Genetics. JRL and RAG declare that Baylor Genetics is a Baylor College of Medicine affiliate that derives revenue from genetic testing. BCM and Miraca Holdings have formed a joint venture with shared ownership and governance of Baylor Genetics which performs clinical microarray analysis and other genomic studies (exome and genome sequencing) for patient and family care. FJS has received research support from Illumina, Pacific Biosciences, and Genentech. JEP is an advisor to MaddieBio. SBM is an advisor to BioMarin, MyOme, and Tenaya Therapeutics. FJS, DEM have received research support and/or consumables from ONT and have received travel funding to speak on behalf of ONT. DEM has received travel support from Pacific Biosciences. DEM is on an advisory board at ONT, a scientific advisory board at Basis Genetics, and holds stock options in both MyOme and Basis Genetics. MJB is the chair of the Scientific Advisory Board of GeneDx and receives funding from the American Society of Human Genetics as the Editor-in-Chief of HGG Advances. JXC receives funding from the American Society of Human Genetics as the Deputy Editor of HGG Advances. DP consults for Ionis Pharmaceuticals. HLR has received rare-disease research funding from Microsoft and Illumina and compensation as a past member of the scientific advisory board of Genome Medical. AODL was a paid consultant to Tome Biosciences, Ono Pharma USA, Addition Therapeutics, Congenica and receives research

funding from Pacific Biosciences. EEE is a scientific advisory board member of Variant Bio, Inc. and is an investigator of the Howard Hughes Medical Institute.

GREGoR Banner Authors:

U01HG011758

Jennifer Posey, Richard Gibbs, James Lupski, Hatoon Al Ali, Elizabeth Atkinson, Sairam Behera, Shaghayegh Beheshti, Eric Boerwinkle, Tugce Bozkurt-Yozgatli, Daniel Calame, Ivan Chinn, Zeynep Coban-Akdemir, Karen Coveler, Zain Dardas, Moez Dawood, Harsha Doddapaneni, Haowei Du, Ruizhi Duan, Iman Egab, Jawid Fatih, Mira Gandhi, Brandon Garcia, Nikhita Gogate, Christopher Grochowski, Jianhong Hu, Minal Jamsandekar, Shalini Jhangiani, Angad Jolly, Parneet Kaur, Ahmed K. Saad, Jesse Levine, Richard Lewis, Yidan Li, Pengfei Liu, Medhat Mahmoud, Dana Marafi, Tadahiro Mitani, Chloe Munderloh, Donna Muzny, Sebastian Ochoa Gonzalez, Piyush Panchal, Shruti Pande, Davut Pehlivan, Archana Rai, Edgar Andres Rivera-Munoz, Aniko Sabo, Evette Scott, Fritz Sedlazeck, V. Reid Sutton, Kim Walker, Lauren Westerfield, Jiaoyang Xu, Bo Yuan, Xinchang Zheng

U01HG011755

Anne O'Donnell-Luria, Heidi Rehm, Michael Talkowski, Siwaar Abouhala, Kaileigh Ahlquist, Mutaz Amin, Christina Austin-Tse, Samantha Baxter, Benjamin Blankenmeister, Philip Boone, Harrison Brand, Colleen Carlston, Celine de Esch, Stephanie DiTroia, Michael Duyzend, Vijay Ganesh, Kiran Garimella, Carmen Glaze, Emily Gropman, Sanna Gudmundsson, Stacey Hall, Yongqing Huang, Julia Klugherz, Katie Larsson, Arthur Lee, Gabrielle Lemire, Jialan Ma, Daniel MacArthur, Brian Mangilog, Daniel Marten, Eva Martinez, Olfa Messaoud, Mariana Moyses, Ashana Neale, Emily O'Heir, Melanie O'Leary, Ikeoluwa Osei-Owusu, Lynn Pais, Alicia Pham, Lindsay Romo, Kathryn Russell, Monica Salani, Kaitlin Samocha, Alba Sanchis-Juan, Jillian Serrano, Gulalai Shah, Moriel Singer-Berk, Mugdha Singh, Hana Snow, Kayla Socarras, Sarah Stenton, Jui-Cheng Tai, Grace VanNoy, Ben Weisburd, Michael Wilson, Monica Wojcik, Isaac Wong, Rachita Yadav

U01HG011762

Stephen Montgomery, Jon Bernstein, Matthew Wheeler, Emily Alsentzer, Raquel Alvarez, Euan Ashley, Themistocles Assimes, Gill Bejerano, Devon Bonner, Denver Bradley, Jennefer Carter, Clarisa Chavez, Ziwei Chen, Salil Deshpande, Sara Emami, Ivy Evergreen, Casey Gifford, Pagé Goddard, John Gorzynski, William Greenleaf, Rodrigo Guarischi-Sousa, Caitlin Harrington, Sohaib Hassan, Tanner Jensen, David Jimenez-Morales, Christopher Jin, Aimee Juan, Jessica Kain, Laura Keehan, Anshul Kundaje, Soumya Kundu, Samuel M. Lancaster, Shruti Marwaha, Dena Matalon, Taylor Maurer, Lauren Meador, Hector Rodrigo Mendez, Alexander Miller, Matthew B. Neu, Thuy-mi P. Nguyen, Jonathan Nguyen, Jeren Olsen, Evin Padhi, Paul Petrowski, Astaria Podesta, Elizabeth Porter, Wanqiong Qiao, Thomas Quertermous, Chloe Reuter, Oriane Rubio, Stuart Scott, Riya Sinha, Kevin S. Smith, Michael Snyder, Brigitte Stark, Suchitra Sudarshan, Christina Tise, Philip Tsao, Rachel Ungar, Isabella Voutos, Juliana Walrod, Ziming Weng, Laurens Wiel, Frank Wong, Yao Yang, Jiye Yu, Jimmy Zhen

U01HG011745

Eric Vilain, Seth Berger, Emmanuèle Délot, Miguel Almalvez, Rishi Aryal, Light Auriga, Rebekah Barrick, Sami Belhadj, Krista Bluske, Leandros Boukas, Andrea J. Cohen, Ya Cui, Ivan de Dios, Meghan Delaney, Jamie Fraser, Vincent Fusaro, John Harting, Megan Hawley, Yun-Hua Hsiao, Amanda Kahn-Kirby, Rachid Karam, Charles Hadley King, Arthur Ko, Wei Li, Bojan Losic, Jonathan LoTempio, Sofia Marmolejos, Robert Nussbaum, Georgia Pitsava, Sarah Savage, Emily Westheimer, Changrui Xiao, Jianhua Zhao

U01HG011744

Michael Bamshad, Chia-Lin Wei, Evan Eichler, Jessica Chong, Kailyn Anderson, Peter Anderson, Sabrina Best, Elizabeth Blue, Kati Buckingham, Silvia Casadei, Yong-Han Cheng, Colleen Davis, Sophia Gibson, William Gordon, Jonas Gustafson, William Harvey, Martha Horike-Pyne, Gail Jarvik, Annelise Mah-Som, Colby Marvin, Francesco Kumara Mastroso, Sean McGee, Heather Mefford, Danny Miller, Miranda Zalusky, Karynne Patterson, Matthew Richardson, Adriana Estela Sedeño-Cortés, Joshua Smith, Olivia Sommerland, Lea Starita, Andrew Stergachis, Elliott Swanson, Jeffrey Weiss, Qian Yi, Christina Zakarian

U24HG011746

Susanne May, Ali Shojaie, Emily Bonkowski, Sarah Conner, Matthew Conomos, Stephanie Gogarten, Ben Heavner, Sarah Nelson, Sheryl Payne, Jaime Prosser, Guanghao Qi, Adrienne Stilp, Catherine Tong, Marsha Wheeler, Quenna Wong

GREGoR Partner Members

Aashish Adhikari, Kinga Bujakowska, Claudia M. B. Carvalho, Ali Crawford, Aimée M. Dudley, Kelly Hagman, Yang I. Li, Jill Moore, Aaron R. Quinlan, Alex Wagner, Bo Xia, S. Stephen Yi

NHGRI

Lisa Chadwick, Christopher Wellington, Sara Currin, Gaby Villard

Writing Group

Moez Dawood, Ben Heavner, Marsha Wheeler, Rachel A. Ungar, Jonathan LoTempio, Laurens Wiel, Claudia M. B. Carvalho, Casey A. Gifford, Susanne May, Danny E. Miller, Heidi L. Rehm, Fritz J. Sedlazeck, Eric Vilain, Anne O'Donnell-Luria, Jennifer E. Posey, Lisa H. Chadwick, Michael J. Bamshad, Stephen B. Montgomery

Supplementary Table 1: Tracking GREGoR Papers With Molecular Diagnoses

Gene	PMID	Functional Work
<i>ACBD6</i>	34582790	No
<i>ACBD6</i>	37951597	Yes
<i>ACCN2</i>	34582790	No
<i>ACOT7</i>	34582790	No
<i>ACTC1</i>	37457373	No
<i>ACTL6A</i>	34582790	No
<i>ACTL6B</i>	39275948	Yes
<i>ACTR1B</i>	39033378	No
<i>ADAM19</i>	34582790	No
<i>ADAMTS15</i>	35962790	Yes
<i>ADSL</i>	34582790	No
<i>AFF3</i>	38811945	Yes
<i>AHDC1</i>	34582790	No
<i>AHDC1</i>	33372375	No
<i>AHDC1</i>	34950897	No
<i>ALS2</i>	34582790	No
<i>AMPD2</i>	34582790	No
<i>ANK3</i>	38988293	No
<i>ANKRD11</i>	34582790	No
<i>AP3B2</i>	34582790	No
<i>AP4B1</i>	34582790	No
<i>APTX</i>	34582790	No
<i>ARAP1</i>	34582790	No
<i>ARFGEF3</i>	38258669	Yes
<i>ARID1B</i>	34582790	No
<i>ARID4A</i>	34582790	No
<i>ARV1</i>	34582790	No
<i>ARX</i>	34582790	No
<i>ASH1L</i>	34582790	No
<i>ASNS</i>	34582790	No
<i>ASPM</i>	34582790	No
<i>ASTN1</i>	34582790	No
<i>ASTN2</i>	34582790	No
<i>ASXL3</i>	34582790	No
<i>ATP1A1</i>	34582790	No
<i>ATP1A3</i>	37043503	Yes
<i>ATP5F1A</i>	34954817	Yes
<i>ATP5F1E</i>	34954817	Yes
<i>ATP5MC3</i>	34954817	Yes
<i>ATP5PO</i>	34954817	Yes
<i>ATP7A</i>	34582790	No
<i>ATRX</i>	34582790	No
<i>BARD1</i>	34582790	No
<i>BHLHA9</i>	36035248	No
<i>BMPER</i>	34582790	No
<i>BRWD3</i>	34582790	No

<i>C2ORF69</i>	34582790	No
<i>CACNA2D2</i>	34582790	No
<i>CAMSAP1</i>	34582790	No
<i>CAPN3</i>	34816580	No
<i>CASP5</i>	37603195	Yes
<i>CBX6</i>	34582790	No
<i>CC2D1B</i>	34582790	No
<i>CCDC39</i>	39606420	No
<i>CCDC40</i>	39606420	No
<i>CCNO</i>	39606420	No
<i>CDK10</i>	34582790	No
<i>CDKL5</i>	35934918	No
<i>CDKL5</i>	34582790	No
<i>CELF2</i>	38258669	Yes
<i>CELSR3</i>	38429302	Yes
<i>CEP290</i>	34582790	No
<i>CEP85L</i>	34582790	No
<i>CFAP46</i>	39606420	No
<i>CFAP47</i>	38633811	Yes
<i>CHASERR</i>	39442041	Yes
<i>CHD2</i>	39442041	Yes
<i>CHD3</i>	34582790	No
<i>CHMP1A</i>	34582790	No
<i>CIT</i>	34582790	No
<i>CLP1</i>	34582790	No
<i>CNTN5</i>	34582790	No
<i>CNTNAP2</i>	34582790	No
<i>COBL</i>	34582790	No
<i>COG3</i>	37711075	Yes
<i>COL6A1</i>	38585825	Yes
<i>COPB1</i>	34582790	No
<i>CREB3</i>	34582790	No
<i>CRYBB2</i>	38990107	No
<i>CTBP1</i>	34582790	No
<i>CTNNA1</i>	38585811	No
<i>CUX1</i>	38585811	No
<i>CYP1B1</i>	38990107	No
<i>DCX</i>	34582790	No
<i>DDC</i>	34582790	No
<i>DDX3X</i>	34582790	No
<i>DEAF1</i>	34582790	No
<i>DENND5A</i>	39174524	Yes
<i>DHCR24</i>	34582790	No
<i>DHX9</i>	37467750	Yes
<i>DLGAP1</i>	34582790	No
<i>DNAAF2</i>	39606420	No
<i>DNAAF4</i>	39606420	No
<i>DNAH1</i>	39606420	No

<i>DNAH1</i>	39606420	No
<i>DNAH11</i>	39606420	No
<i>DNAH5</i>	39606420	No
<i>DNAH6</i>	39606420	No
<i>DNAH6</i>	39606420	No
<i>DNAH8</i>	39606420	No
<i>DNAH9</i>	39606420	No
<i>DNAJC8</i>	34582790	No
<i>DPF2</i>	34582790	No
<i>DRC1</i>	39606420	No
<i>DUSP4</i>	34582790	No
<i>DYRK1A</i>	34582790	No
<i>EEF1A2</i>	34582790	No
<i>ELF4</i>	36477361	Yes
<i>ENPP6</i>	34582790	No
<i>ENTPD1</i>	35471564	Yes
<i>EPG5</i>	34582790	No
<i>EPHA8</i>	34582790	No
<i>ERCC6</i>	34582790	No
<i>ESAM</i>	34582790	No
<i>ESAM</i>	36996813	Yes
<i>EXOSC3</i>	34582790	No
<i>FA2H</i>	34582790	No
<i>FAM120A</i>	34582790	No
<i>FAM91A1</i>	34582790	No
<i>FBP2</i>	38258669	Yes
<i>FBXW11</i>	34582790	No
<i>FER</i>	38585811	No
<i>FGF21</i>	38585811	No
<i>FLNA</i>	39606420	No
<i>FLVCR1</i>	39306721	Yes
<i>FOXG1</i>	34582790	No
<i>FOXI3</i>	37041148	Yes
<i>FOXN4</i>	34582790	No
<i>FRMD7</i>	34582790	No
<i>FSHR</i>	34582790	No
<i>GAS8</i>	39606420	No
<i>GATSL3</i>	34582790	No
<i>GCC2</i>	34582790	No
<i>GCH1</i>	35083481	No
<i>GGPS1</i>	35869884	No
<i>GIN1</i>	34582790	No
<i>GIPR</i>	34582790	No
<i>GIT1</i>	34582790	No
<i>GJC2</i>	34582790	No
<i>GLB1</i>	34582790	No
<i>GLI2</i>	38990107	No
<i>GLI2</i>	34582790	No

<i>GLI3</i>	36035248	No
<i>GMPPB</i>	34582790	No
<i>GNAS</i>	35811283	No
<i>GNAS</i>	38585811	No
<i>GOLGA2</i>	34582790	No
<i>GOLGA4</i>	34582790	No
<i>GPR87</i>	36622818	No
<i>GPT2</i>	34582790	No
<i>GREB1L</i>	37124138	No
<i>GRM7</i>	34582790	No
<i>HECTD3</i>	34582790	No
<i>HECTD4</i>	34582790	No
<i>HECTD4</i>	36401616	No
<i>HEXB</i>	34582790	No
<i>HMGCR</i>	37167966	Yes
<i>HOXD</i>	36035248	No
<i>HOXD13</i>	36035248	No
<i>HPDL</i>	34582790	No
<i>HPS1</i>	34582790	No
<i>HSPB1</i>	33686258	No
<i>HYAL2</i>	34906488	Yes
<i>HYDIN</i>	39606420	No
<i>ITGB8</i>	34582790	No
<i>JRK</i>	34582790	No
<i>KANSL3</i>	38258669	Yes
<i>KCNJ14</i>	34582790	No
<i>KCTD7</i>	34582790	No
<i>KDM2B</i>	34582790	No
<i>KDM5A</i>	34582790	No
<i>KDM5B</i>	39202393	No
<i>KDM5C</i>	34582790	No
<i>KIAA0430</i>	34582790	No
<i>KIF1A</i>	34582790	No
<i>KIF21A</i>	38585811	No
<i>KIF26A</i>	34582790	No
<i>KIF26A</i>	36228617	Yes
<i>KIF5C</i>	38585811	No
<i>KIF7</i>	39606420	No
<i>KIFC3</i>	34582790	No
<i>KLB</i>	38585811	No
<i>L1CAM</i>	34582790	No
<i>LAMA1</i>	34582790	No
<i>LAMB3</i>	34582790	No
<i>LAMC3</i>	34582790	No
<i>LARGE1</i>	34582790	No
<i>LARP7</i>	34582790	No
<i>LCTL</i>	34582790	No
<i>LGI3</i>	35948005	Yes

<i>LPAR6</i>	34582790	No
<i>LRP2</i>	34582790	No
<i>LSS</i>	37157980	No
<i>MAP2K4</i>	38258669	Yes
<i>MAP3K20</i>	38451290	No
<i>MAP3K7</i>	34582790	No
<i>MCM3AP</i>	34582790	No
<i>MCM6</i>	38258669	Yes
<i>MCPH1</i>	34582790	No
<i>MDM1</i>	38868186	Yes
<i>MEGF8</i>	39606420	No
<i>MGAT2</i>	34582790	No
<i>MGP</i>	37923733	Yes
<i>MKS1</i>	34582790	No
<i>MPZ</i>	38585811	No
<i>MRPS25</i>	34582790	No
<i>MRPS25</i>	39606420	No
<i>MRPS27</i>	34582790	No
<i>MTOR</i>	34582790	No
<i>MTSS2</i>	36067766	Yes
<i>MUSK</i>	34816580	No
<i>MYH1</i>	34582790	No
<i>MYH10</i>	38585811	No
<i>NALCN</i>	34582790	No
<i>NANS</i>	34582790	No
<i>NAV2</i>	34582790	No
<i>NES</i>	38585811	No
<i>NETO1</i>	36622818	No
<i>NFE2L3</i>	38258669	Yes
<i>NGEF</i>	34582790	No
<i>NGLY1</i>	34582790	No
<i>NHLRC2</i>	37188825	Yes
<i>NLK</i>	34582790	No
<i>NODAL</i>	38570875	No
<i>NOTCH1</i>	33686258	No
<i>NPHP3</i>	39606420	No
<i>NPR2</i>	36035248	No
<i>NRD1</i>	34582790	No
<i>NSD1</i>	34582790	No
<i>NTNG2</i>	34582790	No
<i>NUAK1</i>	34582790	No
<i>OCLN</i>	34582790	No
<i>OLIG2</i>	38585811	No
<i>OTOA</i>	33492714	No
<i>OTUD6B</i>	35430327	No
<i>PAFAH1B1</i>	34582790	No
<i>PARD3B</i>	34582790	No
<i>PAX5</i>	35094443	No

<i>PCDH18</i>	34582790	No
<i>PDK1L1</i>	39606420	No
<i>PDZD2</i>	34582790	No
<i>PEX6</i>	34582790	No
<i>PGAP3</i>	34582790	No
<i>PHF8</i>	35469323	No
<i>PIK3C2A</i>	34582790	No
<i>PKD1</i>	34582790	No
<i>PLAA</i>	34582790	No
<i>PLCG1</i>	38260438	Yes
<i>PLD3</i>	34582790	No
<i>PLK4</i>	34582790	No
<i>PLXNA1</i>	34582790	No
<i>PNKP</i>	34582790	No
<i>POLR1D</i>	34582790	No
<i>POLR3A</i>	34582790	No
<i>POMGNT1</i>	34582790	No
<i>PPP1R15A</i>	34582790	No
<i>PPP1R21</i>	38356149	No
<i>PPP1R35</i>	34582790	No
<i>PPP1R35</i>	36598158	Yes
<i>PPP1R3F</i>	37531237	Yes
<i>PPP2R1A</i>	34582790	No
<i>PPP2R5C</i>	38258669	Yes
<i>PREX2</i>	34582790	No
<i>PTCHD2</i>	34582790	No
<i>PYCR2</i>	34582790	No
<i>RAC3</i>	35851598	Yes
<i>RAD21</i>	34582790	No
<i>RANBP3L</i>	34582790	No
<i>RARB</i>	37092537	Yes
<i>RASGRF2</i>	34582790	No
<i>RBM10</i>	34582790	No
<i>RCOR3</i>	34582790	No
<i>RNASEH2A</i>	34582790	No
<i>RNASEH2B</i>	34582790	No
<i>RNU4-2</i>	38991538	No
<i>ROBO1</i>	35227688	Yes
<i>ROBO3</i>	38585811	No
<i>ROBO3</i>	34582790	No
<i>RPA1</i>	34582790	No
<i>RSPH4A</i>	39606420	No
<i>RSPH4A</i>	39606420	No
<i>RSPO4</i>	34582790	No
<i>RTN2</i>	34582790	No
<i>RXRA</i>	38990107	No
<i>SCN1A</i>	34582790	No
<i>SCN7A</i>	34582790	No

<i>SEMA3F</i>	38585811	No
<i>SERAC1</i>	34582790	No
<i>SERPINB8</i>	36622818	No
<i>SETX</i>	34582790	No
<i>SF3B1</i>	38258669	Yes
<i>SHANK3</i>	34582790	No
<i>SHROOM4</i>	34582790	No
<i>SHROOOM3</i>	39606420	No
<i>SLC12A5</i>	38585811	No
<i>SLC18A2</i>	34582790	No
<i>SLC19A3</i>	34582790	No
<i>SLC25A45</i>	34582790	No
<i>SLC30A7</i>	35751429	No
<i>SLC30A7</i>	35751429	No
<i>SLC37A1</i>	34582790	No
<i>SLC39A10</i>	34582790	No
<i>SLC4A10</i>	38054405	Yes
<i>SLC5A7</i>	39135055	No
<i>SLC6A1</i>	34582790	No
<i>SLC7A1</i>	34582790	No
<i>SMARCA1</i>	34582790	No
<i>SMC3</i>	38297832	No
<i>SMPD1</i>	34582790	No
<i>SNAPC4</i>	36965478	Yes
<i>SNX14</i>	34582790	No
<i>SORCS2</i>	34582790	No
<i>SOX11</i>	34582790	No
<i>SPAST</i>	36103453	No
<i>SPR</i>	34582790	No
<i>SRD5A3</i>	34582790	No
<i>SRSF1</i>	34582790	No
<i>SSH3</i>	34582790	No
<i>SSTR3</i>	34582790	No
<i>STOML1</i>	34582790	No
<i>STON2</i>	34582790	No
<i>STRTS locus</i>	38714869	Yes
<i>STUB1</i>	34582790	No
<i>SUCLA2</i>	34582790	No
<i>SUFU</i>	34675124	No
<i>SYDE1</i>	36622818	No
<i>SYNE1</i>	38716726	Yes
<i>SYNGAP1</i>	34582790	No
<i>TALDO1</i>	34677006	No
<i>TBC1D23</i>	34582790	No
<i>TBX6</i>	31888956	Yes
<i>TCEAL1</i>	36368327	Yes
<i>TERF2</i>	34582790	No
<i>TGFBR2</i>	38585811	No

<i>THOC6</i>	34582790	No
<i>TLK1</i>	38868186	Yes
<i>TLR7</i>	35477763	Yes
<i>TMOD1</i>	34582790	No
<i>TNRC6B</i>	34582790	No
<i>TOR1A</i>	36757831	No
<i>TPH2</i>	34582790	No
<i>TRAK1</i>	34582790	No
<i>TRAPPC4</i>	34582790	No
<i>TRIM66</i>	34582790	No
<i>TRIT1</i>	34582790	No
<i>TRMT1</i>	34582790	No
<i>TRMT2B</i>	34582790	No
<i>TTC12</i>	39606420	No
<i>TTLL11</i>	34582790	No
<i>TUBA1A</i>	34582790	No
<i>TUBB</i>	38585811	No
<i>TUBB4A</i>	38585811	No
<i>TUBB4A</i>	34582790	No
<i>TUBB6</i>	34582790	No
<i>TUBGCP2</i>	34582790	No
<i>UBE3C</i>	36401616	No
<i>UBR5</i>	38258669	Yes
<i>ULK2</i>	34582790	No
<i>UPS54</i>	34582790	No
<i>VANGL1</i>	38669183	Yes
<i>VANGL2</i>	38669183	Yes
<i>VPS28</i>	34582790	No
<i>VRK3</i>	34582790	No
<i>VSTM2L</i>	34582790	No
<i>WARS2</i>	34582790	No
<i>WDR45B</i>	35322404	No
<i>WDR62</i>	34582790	No
<i>WDR7</i>	34582790	No
<i>WDR73</i>	34582790	No
<i>WDR81</i>	34582790	No
<i>WDR83OS</i>	34582790	No
<i>WNT10B</i>	36035248	No
<i>WRN</i>	35534204	No
<i>WSB1</i>	34582790	No
<i>XRN1</i>	38258669	Yes
<i>YWHAG</i>	38491959	No
<i>ZBTB34</i>	38258669	Yes
<i>ZC4H2</i>	34816580	No
<i>ZFH3</i>	38412861	Yes
<i>ZFH4</i>	39148819	Yes
<i>ZNF462</i>	38585811	No
<i>ZRSR2</i>	34582790	No