

# GeneBreaker: Variant simulation to improve the diagnosis of Mendelian rare genetic diseases

Phillip A. Richmond<sup>1</sup> | Tamar V. Av-Shalom<sup>1</sup> | Oriol Fornes<sup>1</sup> | Bhavi Modi<sup>1</sup> |  
Alison M. Elliott<sup>2</sup> | Wyeth W. Wasserman<sup>1</sup>

<sup>1</sup>Department of Medical Genetics, Center for Molecular Medicine and Therapeutics, BC Children's Hospital Research Institute, University of British Columbia, Vancouver, British Columbia, Canada

<sup>2</sup>Department of Medical Genetics, University of British Columbia, Vancouver, British Columbia, Canada

## Correspondence

Wyeth W. Wasserman, Department of Medical Genetics, Center for Molecular Medicine and Therapeutics, BC Children's Hospital Research Institute, University of British Columbia, Vancouver, BC V5Z 4H4, Canada.

Email: [wyeth@cmmt.ubc.ca](mailto:wyeth@cmmt.ubc.ca)

## Funding information

Michael Smith Foundation for Health Research, Grant/Award Number: 17746; Canadian Institutes of Health Research, Grant/Award Numbers: BOP-149430, PJT-162120; Genome Canada, Grant/Award Number: 255ONT, 275SIL; Genome British Columbia, Grant/Award Numbers: 275SIL, SIP007

## Abstract

Mendelian rare genetic diseases affect 5%–10% of the population, and with over 5300 genes responsible for ~7000 different diseases, they are challenging to diagnose. The use of whole-genome sequencing (WGS) has bolstered the diagnosis rate significantly. The effective use of WGS relies on the ability to identify the disrupted gene responsible for disease phenotypes. This process involves genomic variant calling and prioritization, and is the beneficiary of improvements to sequencing technology, variant calling approaches, and increased capacity to prioritize genomic variants with potential pathogenicity. As analysis pipelines continue to improve, careful testing of their efficacy is paramount. However, real-life cases typically emerge anecdotally, and utilization of clinically sensitive and identifiable data for testing pipeline improvements is regulated and limiting. We identified the need for a gene-based variant simulation framework that can create mock rare disease scenarios, utilizing known pathogenic variants or through the creation of novel gene-disrupting variants. To fill this need, we present GeneBreaker, a tool that creates synthetic rare disease cases with utility for benchmarking variant calling approaches, testing the efficacy of variant prioritization, and as an educational mechanism for training diagnostic practitioners in the expanding field of genomic medicine. GeneBreaker is freely available at <http://GeneBreaker.cmmt.ubc.ca>.

## KEYWORDS

benchmarking, genomics, rare disease, simulation, variant calling, variant interpretation

## 1 | BACKGROUND

Next-generation sequencing, and increasingly third-generation sequencing, has been revolutionary in rare disease diagnosis (Wise et al., 2019). By sequencing the entire genome, millions of variants are identified, prioritized, and then manually curated to

arrive at a diagnosis for the affected individuals. This process occurs both in a familial setting (e.g., trio sequencing of mother–father–proband) as well as in groups of individuals with similar phenotypes (e.g., case series or cohort studies). The analysis process can be broken down into two distinct steps: variant calling and variant interpretation.

Phillip A. Richmond and Tamar V. Av-Shalom contributed equally to this study.

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial License](https://creativecommons.org/licenses/by-nc/4.0/), which permits use, distribution and reproduction in any medium, provided the original work is properly cited and is not used for commercial purposes.

© 2021 The Authors. *Human Mutation* published by Wiley Periodicals LLC

Advances in bioinformatics tools for variant calling have recently expanded the types of variants being detected from single nucleotide variants (SNVs) and small insertions and deletions (indels), to now routinely include copy number variants (CNVs), tandem repeat expansions (REs), mobile element insertions (MEIs), and complex structural variants (SVs) such as inversions, insertions, and translocations. A major challenge within the expansion to these more complex variant types is a large amount of noise and artifacts stemming from limitations of short-reads—with a length of 100–250 base pairs (bp)—attempting to resolve repetitive sequences within the genome. These short reads, even when being sequenced from two ends of a DNA fragment with ~500 bp length, cannot span the repetitive elements in the genome. Specifically, interspersed repeats and short tandem repeats whose length often exceeds 500 bp, cannot be resolved and uniquely mapped against a reference. These mapping ambiguities lead to strict region-based filtering to mitigate noise, in spite of the significant proportion of disease-associated genes overlapping these regions (Ebbert et al., 2019; Goldfeder et al., 2016). As both read length technology and algorithmic approaches continue to evolve, new tools will emerge as candidates to use within diagnostic pipelines (Wenger et al., 2019). While benchmarking the variant calling process in human genomes has been a focus of international consortia, the majority of comparisons focus on evaluations of healthy individuals with well-characterized variant sets (Krusche et al., 2019). In the diagnosis of rare genetic diseases, benchmarks should be focused on assessing the detection capacity of pathogenic or potentially pathogenic variants.

Beyond the landscape of rapidly evolving variant calling approaches is the emergence of a multitude of variant interpretation tools and pipelines. Several *in silico* effect predictors for the assessment of functional variant impact already exist, and research efforts are now adding interpretation capacities for features outside of the coding regions of the genome, for example, splice regulating sequences (Jaganathan et al., 2019). Furthermore, continuously expanding population databases serve as filters to help identify rare genomic events (Karczewski et al., 2020; Lek et al., 2016). Both *in silico* predictors and the population allele frequency of observed variants are used as filters when attempting to identify a pathogenic variant causal for rare genetic disease phenotypes.

Combinations of different variant calling and variant interpretation pipelines are implemented across the world in clinical-grade and research-grade genomic diagnostic laboratories. These approaches utilize tools that consistently receive upgrades to underlying software and databases used within analysis pipelines. As diagnostic pipelines continue to evolve, there is an emerging need for specific performance testing to compare different tools and ensure each new version of an analysis pipeline can identify the disrupted gene in an applied setting. The process of iteratively “spiking” thousands of pathogenic variants into a background set of variants is a common process, especially within research papers publishing novel prioritization methods (e.g., Exomiser; Robinson et al., 2014). However these methods typically draw upon known pathogenic events—usually curated by consortiums such as ClinVar—and are

often limited to SNVs and indels. As multiple classes and genic impacts begin to be explored within the automation space, tools that create unique combinations of rare disease scenarios will become necessary for testing. Combining multiple classes of variants across inheritance patterns is critical to rare disease diagnosis, and the careful creation of such scenarios enables benchmarking “edge” cases within automated pipelines.

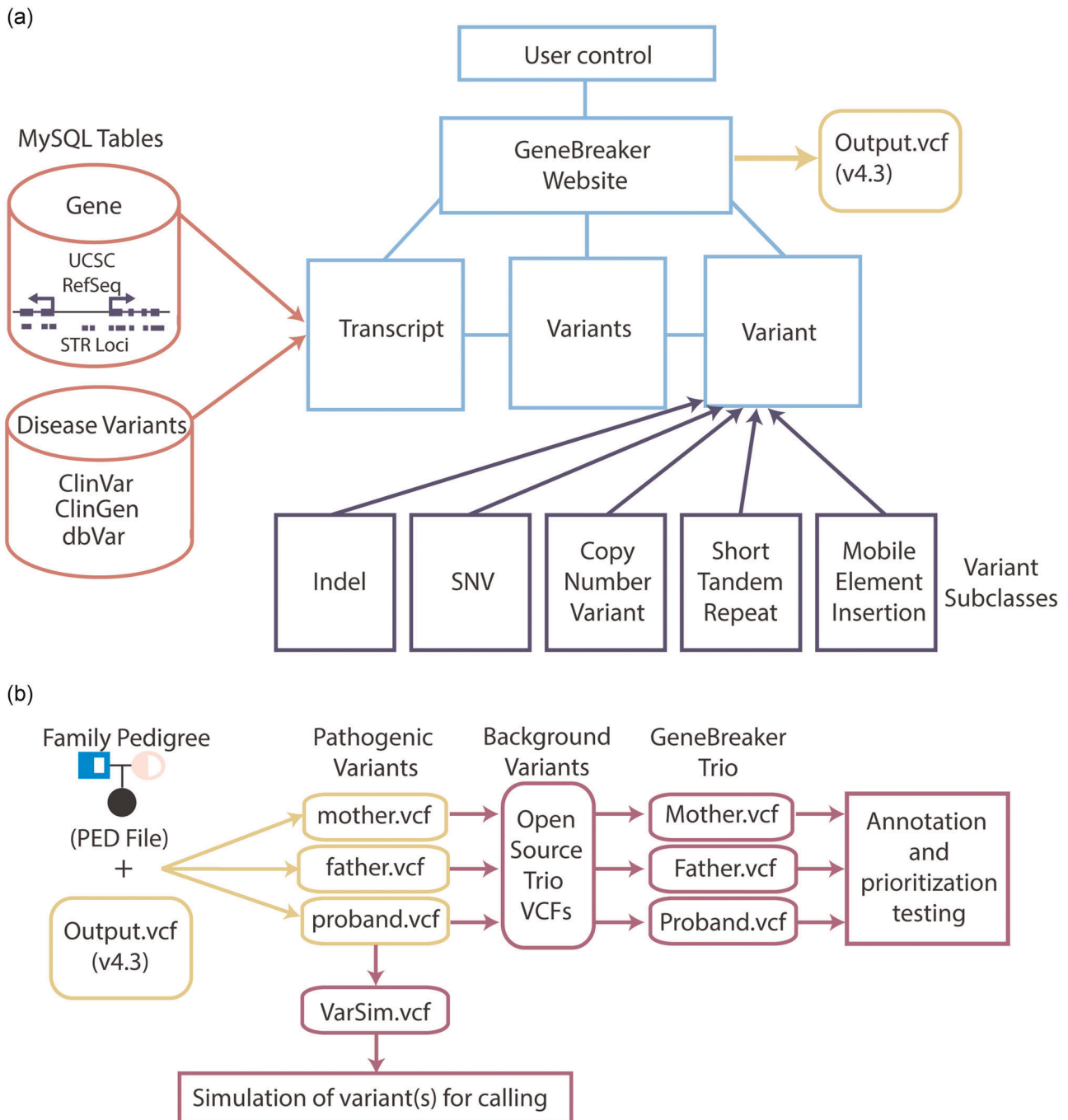
There is a demand for developing the capacity to create synthetic scenarios of rare disease cases for education and training purposes. As genomic medicine advances into the standard of care, there will be a need for easy access to training datasets of rare disease cases of increasing complexity. Institutional policies and guidelines around access and use of sensitive and identifiable data for research purposes beyond that of the specific disease diagnosis vary significantly between different studies and globally (Martani et al., 2019; Raza & Hall, 2017). This means, establishing a universal “standard” of bonafide genomes suitable for personnel training and benchmarking is challenging. Further complicating reanalysis of such data is the possibility of uncovering incidental findings, perhaps affecting either the parents or the proband (Green et al., 2013). In light of incidental findings, this can lead to institutional policies that are strict regarding reanalysis of data post-diagnosis. Even with access to such data for educational purposes, the scale or volume of available data would be limited compared to the potentially infinite possibilities of simulated genomic errors.

To meet these needs and serve the growing community of genomic medicine, we developed GeneBreaker: a simulation tool for Mendelian rare genetic diseases. GeneBreaker has an online web interface for designing custom genetic disease scenarios based on user-guided parameters. It has the capacity to simulate variants of multiple different classes, affecting different genic regions, and can either draw upon known pathogenic variants from resources such as ClinVar (Landrum et al., 2014) and ClinGen (Rehm et al., 2015) or facilitate user creation of novel events. Created variants can be embedded within different familial inheritance models, to model real-life scenarios that may be encountered within clinical or research settings.

## 2 | IMPLEMENTATION

### 2.1 | Architecture

GeneBreaker is a web server deployed as Node.js, which communicates with a REST API that accesses an underlying data repository (storing gene models and known pathogenic variants), and a variant simulation framework written in the Python programming language. User interaction with the online tool guides the stepwise process of variant simulation as follows: (1) select the gene and transcript to be “broken”; (2) simulate the first variant by selecting the variant class and either novel creation or existing pathogenic variant(s); (3) proceed to design the second variant or finish and output a variant call format file (VCF; Figure 1a). Variant creation is done with Mendelian



**FIGURE 1** GeneBreaker overview. (a) Overview of GeneBreaker design framework showing user interaction with the website (light blue), connected MySQL tables (red), underlying variant subclasses (dark blue), and output VCF file (yellow). The user interacts with the GeneBreaker website (light blue) which is connected to hidden components for gene description and variant creation/selection. (b) Downstream benchmarking operations enabled by GeneBreaker including splitting variant amongst VCF files according to user-designed pedigree (yellow), and then either spiking-in the variant within open source trios for annotation and prioritization testing or simulating the proband variant as a full synthetic simulation via VarSim (purple). VCF, variant call format

disease cases in mind, where a single gene/locus is disrupted in either a dominant (single variant) or recessive (one variant per allele) manner. Gene models and known pathogenic variants are stored within a MySQL database, which is used for variant creation. Python code interacts with the MySQL database REST API to extract variant and gene information, and subsequently creates variants based on

user parameters. Code for variant creation and MySQL database interaction can be found here <https://github.com/wassermanlab/GeneBreaker>. While the primary output of the simulator is a VCF file, there is also the capacity to enable downstream benchmarking. Support for downstream benchmarking includes facilitating a transition to a VarSim-compatible VCF file for full synthetic WGS

simulation, and integration into background variant sets for testing annotation and prioritization approaches (Figure 1b). Details about the variant simulation process and downstream benchmarking are below.

## 2.2 | Host web server and underlying data

GeneBreaker is hosted on a virtual web server at the Centre for Molecular Medicine and Therapeutics, with 12 GB of RAM and 4 CPUs running CentOS 7. The variant data within the underlying repository comes from open-source variant catalogs including ClinVar (Landrum et al., 2018), ClinGen (Rehm et al., 2015), and a manually curated set of pathogenic short tandem repeat expansions (Table S1) (Dolzhenko et al., 2020). Gene models come from the RefSeq annotation database provided by the UCSC Genome Browser (NCBI Homo sapiens Annotation Release 109 (March 29, 2018); Haeussler et al., 2019).

## 2.3 | Variant simulation walk-through

The user interaction with the simulation process is stepwise and guided, and a video tutorial detailing the creation process is available at [http://genebreaker.cmmt.ubc.ca/more\\_info](http://genebreaker.cmmt.ubc.ca/more_info). The variant creation process is detailed below.

### 2.3.1 | Initial configuration

Starting at the variant designer page (<http://genebreaker.cmmt.ubc.ca/variants>): (1) user selects reference genome, proband sex, and enters a gene symbol into the “gene” textbox; (2) user clicks “Fetch Transcripts” and all transcripts associated with the gene symbol appear, and the user selects a single transcript which is then displayed in the IGV browser. The primary transcript, as defined by RefSeq Select ([https://www.ncbi.nlm.nih.gov/refseq/refseq\\_select/](https://www.ncbi.nlm.nih.gov/refseq/refseq_select/)), appears with an asterisk to guide transcript selection; and (3) user proceeds to variant creation by clicking “Next.”

### 2.3.2 | Variant creation

From the variant 1 info page: (1) user selects the “Region,” which includes a set of genomic regions: coding, UTR, intronic, and genic (anywhere in the body of the gene). These genomic regions restrict the set of possible variants to those which overlap the defined regions, for example, coding selection means a variant will have to be created over coding regions of the selected gene; (2) user selects the variant “Type”, either choosing from predefined variant sets: ClinVar, ClinGen copy number variant, and short tandem repeat; or novel creation methods: copy number variant, mobile element insertion, indel, single nucleotide variant; and (3) user selects the “Zygosity”: heterozygous or homozygous.

For each of the variant classes, additional information will be required as input. All variant positions are represented in the one-based half-open coordinate system. Users must choose positions for variants that overlap the selected regions. For *ClinVar* and *ClinGen copy number variant*, clicking the “Fetch variants” box will populate the window with coordinate-sorted variants from the ClinVar or ClinGen databases which overlap the defined genomic region. Clicking on a single variant line will select it and enable the user to proceed to the next page. For *copy number variant*, the user specifies the start and end positions, as well as the copy change (deletion or duplication). For *mobile element insertion*, the user specifies the start position and element type: LINE, ALU, or SVA. For *indel*, the user specifies the start position, and the length as a positive integer for insertion of random nucleotides, or as a negative integer for a deletion. For *single nucleotide variant*, there are multiple SNV types that can be selected: stop-loss, missense, nonsense, synonymous, or simply creating alternate alleles from any of the four nucleotides A, T, C, and G. All effects are listed independently of the region selected; however, the effects must comply with the region to proceed. As an example, selecting the intronic region and nonsense variant will give an error. For nonsynonymous variant effects (stop-loss, missense, and nonsense), the variant position must have the capacity to create an amino acid change. As an example, creating a nonsense (premature stop codon) SNV requires that altering the single base at the specified position will change the codon sequence to become TAA, TAG, AGA, or AGG. To facilitate this, the three-frame translation in the IGV window can be displayed by zooming into the nucleotide level, and clicking the gear on the right side of the IGV window to select “Three-frame Translate.” For *short tandem repeat*, clicking on the “Fetch Short tandem repeats” box will populate the window with all short tandem repeats which overlap the genic region. Each STR has the repeat motif, and if the repeat is known to be pathogenic then it displays the number of repeat copies that have to be inserted to be considered damaging. After selecting a repeat, enter the repeat length as an integer value of the number of repeats you want to add or remove, using positive or negative integers, respectively.

After creating a variant, the user can either repeat the process to create a second variant or proceed to Family Info to guide the variant inheritance.

### 2.3.3 | Inheritance modeling

After the creation of variants, a summary page with Family Info appears and a portion of the variant information is displayed including chromosome, position, reference allele sequence, and alternate allele sequence, adhering to specifications from VCF version 4.2 (<https://samtools.github.io/hts-specs/VCFv4.2.pdf>). Below the variant summary is information for each of the individuals in the family which will be included in the output, starting with information from the proband including sex, presence of variant 1 (Var1), presence of variant 2 (Var2), and affected status (check box). The user will then add family members by clicking on the “Add Family Member” box, choosing to add a mother, father, sister, or brother. After adding the

individual, the user will select the boxes for Var1 and Var2 to indicate whether that individual has the variant or not, and their affected status. In the case of homozygous variants, selecting both Var1 and Var2 will add a homozygous variant for that individual. After adding multiple individuals and completing their variant/affected status, the output files can be downloaded by clicking on the "Download Outputs" box. Both the merged VCF containing variant records for all individuals, and the associated pedigree (PED) file can be downloaded.

## 2.4 | Incorporate variants into reads to test detection capacity

Simulation of variants within read sets can be performed in two primary ways: (1) incorporating variants into a reference genome sequence (FASTA format) and then simulating reads from the sequence; or (2) incorporating variants directly into mapped read files. While several tools exist for reference-based incorporation and read simulation, we chose to use VarSim (Mu et al., 2015) due to its ease of use and capacity to simulate other variants alongside the pathogenic variants of interest. Variants were incorporated using VarSim's default configuration with the hg19 (hs37d5) reference genome, background variants from dbSNP common variants version 150 ([https://ftp.ncbi.nih.gov/snp/pre\\_build152/organisms/human\\_9606\\_b150\\_GRCh37p13/](https://ftp.ncbi.nih.gov/snp/pre_build152/organisms/human_9606_b150_GRCh37p13/)), and DGV variants (from VarSim's installation script). Before variant incorporation, VCF files from GeneBreaker may need to be reformatted if they contain a mobile element insertion using the reformatForVarSim.py script (<https://github.com/wassermanlab/GeneBreaker/blob/master/BenchmarkingTransition/FullSimulation/reformatSimToVarSim.py>).

BamSurgeon is a tool for incorporating variants directly into mapped read files (Ewing et al., 2015). GeneBreaker VCF files can be parsed for use within BamSurgeon. We chose to demonstrate the utility of GeneBreaker with VarSim, although both approaches are feasible.

## 2.5 | Spike-in the variant within a trio setting for testing prioritization approaches

Beyond simulating variants for benchmarking variant calling approaches, there is also utility in testing variant prioritization methods. To facilitate this, we collected open source trio PCR-free WGS data from the Polaris Project (<https://github.com/Illumina/Polaris>) and processed it using a standard approach with cutting-edge tools. The data were mapped against the reference genome GRCh37 ([http://www.bcgsc.ca/downloads/genomes/9606/hg19/1000genomes/bwa\\_ind/genome/GRCh37-lite.fa](http://www.bcgsc.ca/downloads/genomes/9606/hg19/1000genomes/bwa_ind/genome/GRCh37-lite.fa)) and GRCh38 ([ftp://ftp.ensembl.org/pub/release-96/fasta/homo\\_sapiens/dna/Homo\\_sapiens.GRCh38.dna\\_sm.primary\\_assembly.fa.gz](ftp://ftp.ensembl.org/pub/release-96/fasta/homo_sapiens/dna/Homo_sapiens.GRCh38.dna_sm.primary_assembly.fa.gz)) using BWA mem (v0.7.17) with default settings (Li & Durbin, 2009). Output SAM files were converted into bam and sorted using Samtools (v1.9; Li et al., 2009). Variant calling was done using DeepVariant (0.10.0; Poplin et al., 2018). Visualizing the mapped reads was done using Integrative Genomics

Viewer (IGV; v2.4.10; Robinson et al., 2011). We applied the same mapping and conversion procedure, using the GRCh37 reference genome, to simulated data from VarSim.

The output is a set of VCF files, one per individual, which serve as background variants for both male and female probands (children) and their parents. These variants are deposited in the online repository alongside other full simulations (see Section "Data Availability Statement"). Combining the background variants with the pathogenic variants from the GeneBreaker tool is managed using bash scripts that match the sex and reference genome. These scripts utilize standard tools including bcftools (Li, 2011), htlib (bgzip and tabix), and a custom reformatting script (<https://github.com/wassermanlab/GeneBreaker/blob/master/BenchmarkingTransition/BuryVariant/reformatSimToDeepVariant.py>). After creating the merged VCF files, we tested them for correct simulation using Exomiser (v12.1.0; Robinson et al., 2014). We searched Exomiser output within each gene-based inheritance table for the known gene using the command line tool grep (e.g., "grep -w ABCD1 -n Exomiser-Output\_AR\_genes.tsv").

## 2.6 | Creation of training scenarios

Hypothetical cases were created and clinical descriptions generated based on clinical experience. Genes and diseases were selected from OMIM for diverse genetic conditions and include rare disorders, including primary immunodeficiencies, inborn errors of metabolism, developmental disorders, and congenital disorders. HPO terms associated with the disease-associated gene were taken from the HPO website (<https://hpo.jax.org/app/>) and selected to include common phenotypes associated with the disease.

## 2.7 | Creation of phenopackets

Phenopackets are an emerging standard, accepted by the Global Alliance for Genomics and Health (GA4GH), for representing phenotypic information in combination with observed variants. For the 10 inheritance cases, as well as the 10 training scenarios, we created phenopacket JSONs using the phenopacket-schema repository (<https://github.com/phenopackets/phenopacket-schema>), following the test example for Python (<https://github.com/phenopackets/phenopacket-schema/tree/master/src/test/python>). The phenopacket JSONs are available on the GeneBreaker website alongside the inheritance testing and training scenarios.

# 3 | RESULTS

## 3.1 | Simulator-created rare disease scenarios

We created rare disease scenarios of varying difficulty to test the efficacy of GeneBreaker simulations and for use within benchmarking scenarios (Table 1). These simulations cover different modes of

**TABLE 1** A set of rare disease scenarios were designed covering a range of Mendelian inheritance patterns and are comprised of either one or two variants affecting a single gene depending on the inheritance pattern

Inheritance	Proband sex	Var1 class	Var1 impact	Source	Var2 class	Var2 impact	Source	Gene	Exomiser rank (noninheritance matched)
Autosomal dominant maternal	Male	SNV	Missense	Published: 28111307	-	-	-	JAK1	2
Autosomal dominant de novo	Male	SNV	Missense	ClinVar: 91028	-	-	-	MSH2	1
Autosomal recessive homozygous	Male	SNV	Missense	Published: 24332264	-	-	-	MALT1	1
Autosomal recessive compound heterozygous	Female	SNV	Missense	Novel	Indel	Frameshift	Novel	CFTR	1
Autosomal recessive compound heterozygous de novo	Male	SNV	Stop-loss	Novel	SNV	Missense	ClinVar: 217654	INPP5E	N/A (7)
X-linked dominant de novo	Female	Indel	Frame-shift	Novel	SNV	Nonsense	Novel	MECP2	1
X-linked recessive homozygous	Male	SNV	Missense	Published: 32202653	-	-	-	WAS	1
X-linked recessive compound heterozygous de novo	Female	Indel	Frame-shift	Novel	SNV	Nonsense	ClinVar: 11696	SLC6A8	N/A (2)
X-linked recessive hemizygous de novo	Male	Indel	Frame-shift	ClinVar: 11303	-	-	-	ABCD1	2
Y-linked de novo	Male	Indel	Coding	Novel	-	-	-	SRY	N/A (N/A)

Note: As many inheritance patterns are sex-specific, the intended proband sex is included. Each variant has the class (SNV or indel), impact on the gene, and source. For variants from the literature or ClinVar, the ClinVar ID or PubMed ID is provided. Exomiser rank is provided, counting the rank of the simulated gene amongst its matching inheritance gene category, and in parentheses is the rank within other inheritance categories.



inheritance, variant classes, and genic impacts. The first set of variants covers several inheritance models for combinations of coding variants, either designed by hand or extracted from ClinVar and other published works. Each of the variants in the table was synthetically generated using the online GeneBreaker interface in either the GRCh37 or GRCh38 genome (Table 1). The variants were then assigned to proband, mother, and father according to the inheritance pattern. Finally, the variants were embedded in the background of open-source trios with matched proband sex (Supporting Information Material). The output from the embedding process is a merged VCF file, which can then be used as input for testing common prioritization workflows, such as the commercial package VarSeq or the open-source Exomiser tool (Robinson et al., 2014).

Beyond creating cases to test inheritance models, we also demonstrate the ability of GeneBreaker to create combinations of variants that are emerging out of anecdotal reports. These variant combinations span multiple classes and genic impacts, and are responsible for the missing heritability in undiagnosed cases (Maroilley & Tarailo-Graovac, 2019). These include a set of cases where pathogenic SNVs and indels lie beyond the coding regions of the gene, and a set of variants which include CNVs, STRs, and MEIs (Table 2).

Finally, we also designed variants within the “dark regions” of the genome, or regions that are inaccessible to standard variant calling pipelines (Ebbert et al., 2019; Goldfeder et al., 2016). We consider these variants important to simulate due to the need to evaluate results from emerging methods capable of genotyping within such regions (Table 2).

### 3.2 | SNV and indel inheritance testing

To the best of our knowledge, there are no currently available open-source tools for prioritizing combinations of different classes of pathogenic variants affecting the same gene. However, the Exomiser tool is a fast and easy-to-use method that can prioritize coding SNVs and indels for Mendelian rare genetic disease cases, and requires as input a merged VCF file, a pedigree (PED) file, and a set of Human Phenotype Ontology (HPO) terms (Robinson et al., 2008, 2014). The HPO terms for a set of selected genes were chosen by matching each gene to a disease using OMIM, and then selecting 4–7 HPO phenotype terms which were common for that disease (<https://hpo.jax.org>; Table S2). We simulated each of the inheritance testing cases (Patients 1–10) and searched the output from Exomiser (Supporting Information Material).

The causal variant was annotated correctly for both user-created SNVs and indels, confirming that our simulation framework for creating novel variants is functional. In 7 out of the 10 cases, the variant was prioritized in the correct inheritance category and was ranked in the top two candidates at the gene level (Table 1). The two scenarios (SLC6A8 and INPP5E) with compound heterozygous de novo inheritance patterns caused issues with Exomiser's interpretation. A compound heterozygous de novo scenario is where a disease-associated allele is inherited from one parent, and a de novo mutation disrupts the other allele of the same gene. In both these scenarios, the

variants were found to be ranked in the top 10 for a dominant inheritance, likely due to the de novo variant taking priority. For instance, the SLC6A8 gene did not show up in the X-linked recessive candidate gene list, but it ranked second in the X-linked dominant gene list. Interestingly, the variant created on the Y chromosome in the SRY gene, which is responsible for 46 XY Sex Reversal 1, was not prioritized by Exomiser. It is unclear at which stage this variant was dropped as a result of Exomiser's inheritance and pathogenicity filtering.

### 3.3 | Testing variant calling

Exomiser is not currently equipped to prioritize CNVs, STRs, and MEIs, and we are not aware of a tool that can integrate these variant classes within inheritance modeling. Therefore, we demonstrate the efficacy of our method by simulating and visualizing a full WGS data set using the larger variants and the set of variants within the dark regions of the genome (Table 2). Using VarSim, the 10 variants from the CNV/MEI/STR and dark genome categories were simulated in a single VarSim run. Each of the regions where variants were integrated was visualized with the IGV (Robinson et al., 2011) to validate the variant was simulated correctly at the read level (Figures S1 and S2). An example of the heterozygous duplication of part of the DMD gene shows the expected increase in read coverage over the simulated variant (Figure S1a), and the LINE1 transposable element insertion within the intron of SLC17A5 has the expected signal of soft-clipped reads both upstream and downstream of the insertion site (Figure S1b). For the variants in the dark genome, a homozygous deletion in SMN2 can be visualized, even though the observed reads are not mapping uniquely to the region (Figure S2a). Finally, a four base-pair coding deletion within CFC1 appears in the ambiguously mapped reads, confirming previous reports that specialized methods may be able to locate deletions within these dark and camouflaged regions (Ebbert et al., 2019; Figure S2b).

### 3.4 | Training scenarios

To emphasize the capacity for GeneBreaker-created scenarios to be utilized within training the next generation of genome medicine practitioners, we created 10 hypothetical genetic disorder scenarios. Each hypothetical case has a causal gene, causal variants drawn from ClinVar, HPO terms, and patient descriptions including relevant family history and clinical findings (Table 3; Supporting Information Material). For these 10 cases, they can be utilized within educational materials that focus on variant interpretation in rare disease diagnosis.

## 4 | DISCUSSION

The diagnosis of rare genetic diseases will continue to improve as novel methods for calling, interpreting, and prioritizing variants emerge and become deployed in a diagnostic setting. Many of the

**TABLE 2** Simulated scenarios for noncoding SNVs and indels; CNVs, MEIs, and STRs; and variants in the dark genome

Inheritance	Proband sex	Var1 class	Var1 impact	Source	Var2 class	Var2 impact	Source	Gene
<i>Noncoding SNVs and indels</i>								
Autosomal recessive compound heterozygous	Male	SNV	Intronic	ClinVar: 99243	Indel	Frameshift	Novel	ABCA4
Autosomal recessive compound heterozygous	Female	SNV	Intronic	ClinVar: 6770	SNV	Nonsense	Novel	PEX10
Autosomal dominant de novo	Male	Indel	Intronic	ClinVar: 411336	-	-	-	APC
Autosomal dominant de novo	Female	SNV	UTR	ClinVar: 556698	-	-	-	MMACHC
Autosomal recessive compound heterozygous	Male	SNV	UTR	ClinVar: 412264	SNV	Nonsense	ClinVar: 523079	BBS9
<i>CNVs, MEIs, and STRs</i>								
Autosomal dominant de novo	Female	CNV (partial DEL)	Coding	Novel	-	-	-	SLC2A1
Autosomal recessive compound heterozygous	Female	CNV (full DEL)	Coding	Novel	Indel	Frameshift	Novel	HEXA
Autosomal recessive homozygous	Female	MEI	Intron	Published: 28187749	-	-	-	SLC17A5
X-linked dominant de novo	Female	MEI	Coding	Novel	-	-	-	IKBKG
Autosomal recessive homozygous	Female	STR	UTR	Published: 30970188	-	-	-	GLS
X-linked recessive compound heterozygous	Female	CNV	Coding	Novel	SNV	Nonsense	ClinVar: 282841	DMD
<i>Variants in the dark genome</i>								
Autosomal recessive homozygous	Female	CNV (partial DEL)	Coding	Published: 32066871	-	-	-	SMN2
Autosomal dominant de novo	Female	Indel	Frameshift	Novel	-	-	-	CFC1
Autosomal dominant de novo	Female	Indel	Frameshift	Novel	-	-	-	MAF
Autosomal dominant de novo	Female	Indel	Frameshift	Novel	MEI	Coding	Novel	RPGR

Note: The inheritance, sex, and variant information for variant 1 and variant 2 are shown, including their class, impact, and source (novel, ClinVar ID, or PubMed ID).



**TABLE 3** Rare disease scenarios created for teaching purposes, labeled as case 1–10, with the sex (M = male, F = female), inheritance model, gene and disease association, ClinVar ID for variant 1, ClinVar ID for variant 2 (if applicable, otherwise “-”), HPO IDs, and HPO phenotype terms

Case No.	Sex	Inheritance	Gene; disease	Var 1		Var 2		Features (matched to HPO terms)
				ClinVar ID	ClinVar ID	ClinVar ID	HPO terms	
1	M	X-linked recessive hemizygous	WAS; Wiskott-Aldrich syndrome	424355	-	"HP:0000964," "HP:0001873," "HP:0005537," "HP:0000388," and "HP:0002573"	-	Eczema, thrombocytopenia, small platelets, otitis media, and hematochezia. The family history is significant for a maternal uncle with autoimmune problems and a recent diagnosis of lymphoma.
2	M	Autosomal dominant de novo	MSH2; Lynch syndrome	91055	-	"HP:0040275," "HP:0012378," and "HP:0002027"	-	Adenocarcinoma of the large intestine, fatigue, and abdominal pain. Negative family history.
3	F	Autosomal recessive homozygous	MALT1; immunodeficiency 12	662739	-	"HP:0000964," "HP:0001581," "HP:0004386," "HP:0002090," and "HP:0002205"	-	Eczema, recurrent skin infections, gastrointestinal inflammation, pneumonia, and recurrent respiratory infections. Parents are consanguineous (first cousins).
4	M	Autosomal recessive compound heterozygous	CFTR; cystic fibrosis	618928	554293	"HP:0002024," "HP:0002206," "HP:0002205," and "HP:0001738"	-	Malabsorption, pulmonary fibrosis, recurrent respiratory infections, and exocrine pancreatic insufficiency. Parents are unrelated and of British ancestry.
5	F	X-linked dominant de novo	MECP2; Rett syndrome	424171	-	"HP:0001250," "HP:0001257," "HP:0005484," and "HP:0002376"	-	Seizure, spasticity, postnatal microcephaly, and developmental regression. Negative family history.
6	M	X-linked recessive hemizygous de novo	ABCD1; adrenoleukodystrophy	458629	-	"HP:0001250," "HP:0000709," "HP:0008207," "HP:0002180," and "HP:0002500"	-	Seizure, psychosis, primary adrenal insufficiency, neurodegeneration, and abnormality of the cerebral white matter. Negative family history.
7	F	Autosomal dominant paternal	EP300; Rubinstein Taybi syndrome Type 2	666310	-	"HP:002342," "HP:0002553," "HP:0000494," "HP:0000448," "HP:0000347," and "HP:0011304"	-	Moderate intellectual disability, arched eyebrows, downslanting palpebral fissures, prominent nose, micrognathia, and broad thumbs. Father with significant learning problems and similar facial features.
8	M	Autosomal dominant de novo	PTPN11; Noonan syndrome	40488	-	"HP:0004322," "HP:0032318," "HP:0000475," "HP:0000368," "HP:0000316," "HP:0000445," and "HP:0003196"	-	Short stature, congenital heart disease, wide neck, low set posteriorly rotated ears, widely spaced eyes, and short and broad nose. Negative family history.

TABLE 3 (Continued)

Case No.	Sex	Inheritance	Gene; disease	Var 1 ClinVar ID	Var 2 ClinVar ID	HPO terms	Features (matched to HPO terms)
9	F	Autosomal dominant paternal	COL2A1; Stickler syndrome	449397	-	"HP:0011800," "HP:0008625," "HP:0000545," and "HP:0000175"	Flat midface, sensorineural hearing loss, myopia, and cleft palate. Father with severe myopia.
10	M	Autosomal dominant maternal	ANKRD11; KBG syndrome	658274	-	"HP:0001249," "HP:0004322," "HP:0000252," "HP:0000325," "HP:0000400," "HP:0000316," "HP:0000637," and "HP:0001572"	Intellectual disability, short stature, microcephaly, triangular face, large and prominent ears, hypertelorism, long palpebral fissures, and macrodontia. Mother similarly affected.

previously challenging variants to genotype, due to variant complexity or existing within “dark” genomic regions, are now regularly identified in WGS datasets. Consequently, there is a critical need for testing the improved pipelines to ensure that such improvements are implemented correctly. While real patient data is paramount for testing the efficacy of a pipeline, access to such data is often prohibitive due to data use restrictions and is limited by the number of observed cases with available data. With simulation, infinite combinations of any possible genomic variant(s) can be designed, from nonsense SNVs, to intronic MEIs, and every combination in between.

The introduction of GeneBreaker addresses an unmet need for simulated rare disease cases in a broadly accessible format. GeneBreaker is deployed as a free-to-use website, enabling user creation of pathogenic variants. Downstream of variant simulation, the tool also supports the transition into benchmarking either variant interpretation or variant calling analysis pipelines. We tested the efficacy of GeneBreaker and the downstream benchmarking transition by simulating rare disease scenarios covering different inheritance models, variant classes, genomic regions, and genic impacts. Using Exomiser, we validated that the variants we simulate have the expected impact, and exposed some limitations in the ability to correctly prioritize variants with challenging inheritance patterns, such as the compound heterozygous de novo pattern. Our example using Exomiser highlights that even for combinations of SNVs and indels, inheritance testing can still be improved upon. Larger, more complex variants were visualized in IGV to ensure their correct simulation, confirming that VarSim is a viable option for whole data set simulation to test variant calling capacity. All of the simulated cases are available online ([http://genebreaker.cmm.ubc.ca/premade\\_cases](http://genebreaker.cmm.ubc.ca/premade_cases)) and can serve as a starting point for benchmarking.

GeneBreaker is not intended as another genome sequencing data simulator. Such simulation has been broadly explored over the past decade, with a rich collection of tools available (Escalona et al., 2016). The narrow scope of GeneBreaker is placed upon the creation of rare disease simulated genomes, which is achieved by focusing on the generation of diverse forms of genetic disruptions, which can be embedded within real or simulated genome sequencing data. This focus has particular value for two use cases: careful edge-testing of analysis pipelines and training of interpretation specialists.

When it comes to changing software within a clinical diagnostic pipeline—even if it only involves upgrading to a newer version of an existing package—there must be rigorous testing to ensure that the modifications do not break the pipeline. Any modification to a standard operating procedure must be tested to ensure that it is still capable of performing at or above the existing diagnostic capacity. With the increased adoption of the reference genome version GRCh38, many pipelines currently utilizing older reference genomes will need mechanisms to test their correct functionality in GRCh38 before the transition. Adopting a new reference genome can sometimes have unanticipated side-effects, as was highlighted with an analysis of missing variant calls from realigning WGS datasets in the UK Biobank (Jia et al., 2020). A diverse set of variant scenarios

involving several inheritance models, genic impacts, and mix of novel and known pathogenic variants has utility for testing these continuously updated pipelines, without the challenge of handling sensitive patient data. GeneBreaker is designed with the goal of creating carefully thought-out edge testing cases and is not built as a genotype-to-phenotype prioritization benchmarking tool. The process of testing gene-to-phenotype associations and rankings is better handled by considering only known pathogenic events and utilizing thousands of simulations. The careful creation of complex scenarios is the focus of the benchmarking aspect of GeneBreaker.

GeneBreaker has value beyond benchmarking as a resource for training a new generation of genome analysts. Rare genetic diseases affect a sizable portion of the population, and as WGS moves into the standard of care, many medical professionals will need hands-on training in the utilization of this technology. Having synthetic cases, either at the merged variant set or raw data levels, is imperative. We encourage those developing educational materials for the analysis of rare disease genomes to consider using the simulation capacity of GeneBreaker as a training tool. To emphasize the teaching aspect of GeneBreaker, and to allow rapid adoption into the educational setting, we created 10 hypothetical genetic disorder scenarios, complete with phenotypes, patient descriptions, family history, and merged variant sets, all available online ([genebreaker.cmmt.ubc.ca](http://genebreaker.cmmt.ubc.ca)). We envision these materials to be invaluable in training individuals working at many institutions, both academic and commercial, who are establishing genome sequencing protocols for rare disease diagnosis.

Future work on GeneBreaker will focus on expanding the simulation capacity to include additional complex variant classes (e.g., inversions and translocations) and variants beyond the genic regions known to disrupt gene regulation. Examples of disruptions to regulatory elements include mutations affecting chromatin organization and enhancers (Lupiáñez et al., 2016; Smith & Shilatifard, 2014). A major challenge in extending to regulatory elements is that the relevant genomic regions critical for the regulation of a gene are difficult to define. As these genome annotations improve, we look forward to integrating them into GeneBreaker.

We hope that GeneBreaker is adopted by the growing community of researchers and clinicians who are utilizing WGS in the diagnosis of rare genetic diseases. Feedback is appreciated as we continue to improve the software and simulation capacities. GeneBreaker is available for exploration at <http://GeneBreaker.cmmt.ubc.ca>.

## ACKNOWLEDGMENTS

The authors would like to acknowledge computational infrastructure support from the Compute Canada and the University of British Columbia Advanced Research Computing organizations. The authors also thank members of the Wasserman lab for testing the software. Phillip A. Richmond, Tamar V. Av-Shalom, Bhavi Modi, and Wyeth W. Wasserman were supported by grants from the Canadian Institutes of Health Research (BOP-149430 and

PJT-162120); the Genome Canada and Genome British Columbia (255ONT and 275SIL); the Michael Smith Foundation for Health Research (17746); and the Genome British Columbia (SIP007).

## CONFLICT OF INTERESTS

The authors declare that there are no conflict of interests.

## AUTHOR CONTRIBUTIONS

Phillip A. Richmond, Tamar V. Av-Shalom, and Wyeth W. Wasserman devised the project. Tamar V. Av-Shalom and Oriol Fornes coded the database schema for gene and variant extraction. Tamar V. Av-Shalom developed the variant creation software and GeneBreaker website. Bhavi Modi and Phillip A. Richmond curated variants for simulation. Phillip A. Richmond generated synthetic datasets, analyzed the data, and developed downstream benchmarking capacity. Alison M. Elliott developed hypothetical case scenarios. Phillip A. Richmond, Tamar V. Av-Shalom, and Wyeth W. Wasserman wrote the manuscript, and all authors approved the manuscript.

## DATA AVAILABILITY STATEMENT

The code for the variant simulation, website deployment, and downstream benchmarking can be found on GitHub: <https://github.com/wassermanlab/GeneBreaker>. The GeneBreaker website can be accessed at <http://genebreaker.cmmt.ubc.ca>. Premade variant sets can be found on the GeneBreaker website and are additionally hosted at Zenodo: <https://doi.org/10.5281/zenodo.3829960>. Training cases can also be found on the GeneBreaker website and are additionally hosted at Zenodo.

## WEB RESOURCES

- NCBI RefSeq Select: [https://www.ncbi.nlm.nih.gov/refseq/refseq\\_select/](https://www.ncbi.nlm.nih.gov/refseq/refseq_select/)
- GRCh37 Reference Genome: [http://www.bcgsc.ca/downloads/genomes/9606/hg19/1000genomes/bwa\\_ind/genome/GRCh37-lite.fa](http://www.bcgsc.ca/downloads/genomes/9606/hg19/1000genomes/bwa_ind/genome/GRCh37-lite.fa)
- GRCh38 Reference Genome: [ftp://ftp.ensembl.org/pub/release-96/fasta/homo\\_sapiens/dna/Homo\\_sapiens.GRCh38.dna\\_sm.primary\\_assembly.fa.gz](ftp://ftp.ensembl.org/pub/release-96/fasta/homo_sapiens/dna/Homo_sapiens.GRCh38.dna_sm.primary_assembly.fa.gz)

## REFERENCES

- Dolzhenko, E., Bennett, M. F., Richmond, P. A., Trost, B., Chen, S., van Vugt, J. J. F. A., Nguyen, C., Narzisi, G., Gainullin, V. G., Gross, A. M., Lajoie, B. R., Taft, R. J., Wasserman, W. W., Scherer, S. W., Veldink, J. H., Bentley, D. R., Yuen, R. K. C., Bahlo, M., & Eberle, M. A. (2020). ExpansionHunter Denovo: A computational method for locating known and novel repeat expansions in short-read sequencing data. *Genome Biology*, 21(1), 102. <https://doi.org/10.1186/s13059-020-02017-z>
- Ebbert, M. T. W., Jensen, T. D., Jansen-West, K., Sens, J. P., Reddy, J. S., Ridge, P. G., Kauwe, J. S. K., Belzil, V., Pregent, L., Carrasquillo, M. M., Keene, D., Larson, E., Crane, P., Asmann, Y. W., Ertekin-Taner, N., Younkin, S. G., Ross, O. A., Rademakers, R., Petrucelli, L., & Fryer, J. D. (2019). Systematic analysis of dark and camouflaged

- genes reveals disease-relevant genes hiding in plain sight. *Genome Biology*, 20(1), 97. <https://doi.org/10.1186/s13059-019-1707-2>
- Escalona, M., Rocha, S., & Posada, D. (2016). A comparison of tools for the simulation of genomic next-generation sequencing data. *Nature Reviews Genetics*, 17(8), 459–469. <https://doi.org/10.1038/nrg.2016.57>
- Ewing, A. D., Houlahan, K. E., Hu, Y., Ellrott, K., Caloian, C., Yamaguchi, T. N., Bare, J. C., P'ng, C., Waggott, D., Sabelnykova, V. Y., Kellen, M. R., Norman, T. C., Haussler, D., Friend, S. H., Stolovitzky, G., Margolin, A. A., Stuart, J. M., & Boutros, P. C. (2015). Combining tumor genome simulation with crowdsourcing to benchmark somatic single-nucleotide-variant detection. *Nature Methods*, 12(7), 623–630. <https://doi.org/10.1038/nmeth.3407>
- Goldfeder, R. L., Priest, J. R., Zook, J. M., Grove, M. E., Waggott, D., Wheeler, M. T., Salit, M., & Ashley, E. A. (2016). Medical implications of technical accuracy in genome sequencing. *Genome Medicine*, 8(1), 24. <https://doi.org/10.1186/s13073-016-0269-0>
- Green, R. C., Berg, J. S., Grody, W. W., Kalia, S. S., Korf, B. R., Martin, C. L., McGuire, A. L., Nussbaum, R. L., O'Daniel, J. M., Ormond, K. E., Rehm, H. L., Watson, M. S., Williams, M. S., & Biesecker, L. G. (2013). ACMG recommendations for reporting of incidental findings in clinical exome and genome sequencing. *Genetics in Medicine*, 15(7), 565–574. <https://doi.org/10.1038/gim.2013.73>
- Haeussler, M., Zweig, A. S., Tyner, C., Speir, M. L., Rosenbloom, K. R., Raney, B. J., Lee, C. M., Lee, B. T., Hinrichs, A. S., Gonzalez, J. N., Gibson, D., Diekhans, M., Clawson, H., Casper, J., Barber, G. P., Haussler, D., Kuhn, R. M., & Kent, W. J. (2019). The UCSC Genome Browser database: 2019 update. *Nucleic Acids Research*, 47(D1), D853–D858. <https://doi.org/10.1093/nar/gky1095>
- Jaganathan, K., Kyriazopoulou Panagiotopoulou, S., McRae, J. F., Darbandi, S. F., Knowles, D., Li, Y. I., Kosmicki, J. A., Arbelaez, J., Cui, W., Schwartz, G. B., Chow, E. D., Kanterakis, E., Gao, H., Kia, A., Batzoglou, S., Sanders, S. J., & Farh, K. K. H. (2019). Predicting splicing from primary sequence with deep learning. *Cell*, 176(3), 535–548. <https://doi.org/10.1016/j.cell.2018.12.015>
- Jia, T., Munson, B., Lango Allen, H., Ideker, T., & Majithia, A. R. (2020). Thousands of missing variants in the UK Biobank are recoverable by genome realignment. *Annals of Human Genetics*, 84(3), 214–220. <https://doi.org/10.1111/ahg.12383>
- Karczewski, K. J., Francioli, L. C., Tiao, G., Cummings, B. B., Alföldi, J., Wang, Q., Collins, R. L., Laricchia, K. M., Ganna, A., Birnbaum, D. P., Gauthier, L. D., Brand, H., Solomonson, M., Watts, N. A., Rhodes, D., Singer-Berk, M., England, E. M., Seaby, E. G., Kosmicki, J. A., ... MacArthur, D. G. (2020). The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature*, 581(7809), 434–443. <https://doi.org/10.1038/s41586-020-2308-7>
- Krusche, P., Trigg, L., Boutros, P. C., Mason, C. E., De La Vega, F. M., Moore, B. L., Gonzalez-Porta, M., Eberle, M. A., Tezak, Z., Lababidi, S., Truty, R., Asimenos, G., Funke, B., Fleharty, M., Chapman, B. A., Salit, M., & Zook, J. M. (2019). Best practices for benchmarking germline small-variant calls in human genomes. *Nature Biotechnology*, 37(5), 555–560. <https://doi.org/10.1038/s41587-019-0054-x>
- Landrum, M. J., Lee, J. M., Benson, M., Brown, G. R., Chao, C., Chitipiralla, S., Gu, B., Hart, J., Hoffman, D., Jang, W., Karapetyan, K., Katz, K., Liu, C., Maddipatla, Z., Malheiro, A., McDaniel, K., Ovetsky, M., Riley, G., Zhou, G., ... Maglott, D. R. (2018). ClinVar: Improving access to variant interpretations and supporting evidence. *Nucleic Acids Research*, 46(D1), D1062–D1067. <https://doi.org/10.1093/nar/gkx1153>
- Landrum, M. J., Lee, J. M., Riley, G. R., Jang, W., Rubinstein, W. S., Church, D. M., & Maglott, D. R. (2014). ClinVar: Public archive of relationships among sequence variation and human phenotype. *Nucleic Acids Research*, 42(Database issue), D980–D985. <https://doi.org/10.1093/nar/gkt1113>
- Lek, M., Karczewski, K. J., Minikel, E. V., Samocha, K. E., Banks, E., Fennell, T., O'donnell-Luria, A. H., Ware, J. S., Hill, A. J., Cummings, B. B., Tukiainen, T., Birnbaum, D. P., Kosmicki, J. A., Duncan, L. E., Estrada, K., Zhao, F., Zou, J., Pierce-Hoffman, E., Berghout, J., ... Exome Aggregation Consortium (2016). Analysis of protein-coding genetic variation in 60,706 humans. *Nature*, 536(7616), 285–291. <https://doi.org/10.1038/nature19057>
- Li, H. (2011). A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics*, 27(21), 2987–2993. <https://doi.org/10.1093/bioinformatics/btr509>
- Li, H., & Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*, 25(14), 1754–1760. <https://doi.org/10.1093/bioinformatics/btp324>
- Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., & Durbin, R. (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, 25(16), 2078–2079. <https://doi.org/10.1093/bioinformatics/btp352>
- Lupiáñez, D. G., Spielmann, M., & Mundlos, S. (2016). Breaking TADs: How alterations of chromatin domains result in disease. *Trends in Genetics*, 32(4), 225–237. <https://doi.org/10.1016/j.tig.2016.01.003>
- Maroille, T., & Tarailo-Graovac, M. (2019). Uncovering missing heritability in rare diseases. *Genes*, 10(4), 275. <https://doi.org/10.3390/genes10040275>
- Martani, A., Geneviève, L. D., Pauli-Magnus, C., McLennan, S., & Elger, B. S. (2019). Regulating the secondary use of data for research: Arguments against genetic exceptionalism. *Frontiers in Genetics*, 10, 1254. <https://doi.org/10.3389/fgene.2019.01254>
- Mu, J. C., Mohiyuddin, M., Li, J., Bani Asadi, N., Gerstein, M. B., Abyzov, A., Wong, W. H., & Lam, H. Y. K. (2015). VarSim: A high-fidelity simulation and validation framework for high-throughput genome sequencing with cancer applications. *Bioinformatics*, 31(9), 1469–1471. <https://doi.org/10.1093/bioinformatics/btu828>
- Poplin, R., Chang, P.-C., Alexander, D., Schwartz, S., Colthurst, T., Ku, A., Newburger, D., Dijamco, J., Nguyen, N., Afshar, P. T., Gross, S. S., Dorfman, L., McLean, C. Y., & DePristo, M. A. (2018). A universal SNP and small-indel variant caller using deep neural networks. *Nature Biotechnology*, 36(10), 983–987. <https://doi.org/10.1038/nbt.4235>
- Raza, S., & Hall, A. (2017). Genomic medicine and data sharing. *British Medical Bulletin*, 123(1), 35–45. <https://doi.org/10.1093/bmb/ldx024>
- Rehm, H. L., Berg, J. S., Brooks, L. D., Bustamante, C. D., Evans, J. P., Landrum, M. J., Ledbetter, D. H., Maglott, D. R., Martin, C. L., Nussbaum, R. L., Plon, S. E., Ramos, E. M., Sherry, S. T., & Watson, M. S. (2015). ClinGen – The Clinical Genome Resource. *New England Journal of Medicine*, 372(23), 2235–2242. <https://doi.org/10.1056/nejmsr1406261>
- Robinson, P. N., Köhler, S., Bauer, S., Seelow, D., Horn, D., & Mundlos, S. (2008). The Human Phenotype Ontology: A tool for annotating and analyzing human hereditary disease. *American Journal of Human Genetics*, 83(5), 610–615. <https://doi.org/10.1016/j.ajhg.2008.09.017>
- Robinson, P. N., Köhler, S., Oellrich, A., Wang, K., Mungall, C. J., Lewis, S. E., Washington, N., Bauer, S., Seelow, D., Krawitz, P., Gilissen, C., Haendel, M., & Smedley, D. (2014). Improved exome prioritization of disease genes through cross-species phenotype comparison. *Genome Research*, 24(2), 340–348. <https://doi.org/10.1101/gr.160325.113>
- Robinson, J. T., Thorvaldsdóttir, H., Winckler, W., Guttman, M., Lander, E. S., Getz, G., & Mesirov, J. P. (2011). Integrative genomics viewer. *Nature Biotechnology*, 29(1), 24–26. <https://doi.org/10.1038/nbt.1754>

- Smith, E., & Shilatifard, A. (2014). Enhancer biology and enhanceropathies. *Nature Structural & Molecular Biology*, 21(3), 210–219. <https://doi.org/10.1038/nsmb.2784>
- Wenger, A. M., Peluso, P., Rowell, W. J., Chang, P.-C., Hall, R. J., Concepcion, G. T., Ebler, J., Fungtammasan, A., Kolesnikov, A., Olson, N. D., Töpfer, A., Alonge, M., Mahmoud, M., Qian, Y., Chin, C. S., Phillippy, A. M., Schatz, M. C., Myers, G., DePristo, M. A., ... Hunkapiller, M. W. (2019). Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome. *Nature Biotechnology*, 37(10), 1155–1162. <https://doi.org/10.1038/s41587-019-0217-9>
- Wise, A. L., Manolio, T. A., Mensah, G. A., Peterson, J. F., Roden, D. M., Tamburro, C., Williams, M. S., & Green, E. D. (2019). Genomic medicine for undiagnosed diseases. *Lancet*, 394(10197), 533–540. [https://doi.org/10.1016/S0140-6736\(19\)31274-7](https://doi.org/10.1016/S0140-6736(19)31274-7)

## SUPPORTING INFORMATION

Additional Supporting Information may be found online in the supporting information tab for this article.

**How to cite this article:** Richmond PA, Av-Shalom TV, Fornes O, Modi B, Elliott AM, Wasserman WW. GeneBreaker: Variant simulation to improve the diagnosis of Mendelian rare genetic diseases. *Human Mutation*. 2021;42:346–358. <https://doi.org/10.1002/humu.24163>