
StrVCTVRE: A supervised learning method to predict the pathogenicity of human genome structural variants

Authors

Andrew G. Sharo, Zhiqiang Hu,
Shamil R. Sunyaev, Steven E. Brenner

Correspondence

brenner@compbio.berkeley.edu (S.E.B.),
sharo@compbio.berkeley.edu (A.G.S.)



StrVCTVRE: A supervised learning method to predict the pathogenicity of human genome structural variants

Andrew G. Sharo,^{1,2,*} Zhiqiang Hu,^{2,3} Shamil R. Sunyaev,^{4,5} and Steven E. Brenner^{1,2,3,*}

Summary

Whole-genome sequencing resolves many clinical cases where standard diagnostic methods have failed. However, at least half of these cases remain unresolved after whole-genome sequencing. Structural variants (SVs; genomic variants larger than 50 base pairs) of uncertain significance are the genetic cause of a portion of these unresolved cases. As sequencing methods using long or linked reads become more accessible and SV detection algorithms improve, clinicians and researchers are gaining access to thousands of reliable SVs of unknown disease relevance. Methods to predict the pathogenicity of these SVs are required to realize the full diagnostic potential of long-read sequencing. To address this emerging need, we developed StrVCTVRE to distinguish pathogenic SVs from benign SVs that overlap exons. In a random forest classifier, we integrated features that capture gene importance, coding region, conservation, expression, and exon structure. We found that features such as expression and conservation are important but are absent from SV classification guidelines. We leveraged multiple resources to construct a size-matched training set of rare, putatively benign and pathogenic SVs. StrVCTVRE performs accurately across a wide SV size range on independent test sets, which will allow clinicians and researchers to eliminate about half of SVs from consideration while retaining a 90% sensitivity. We anticipate clinicians and researchers will use StrVCTVRE to prioritize SVs in probands where no SV is immediately compelling, empowering deeper investigation into novel SVs to resolve cases and understand new mechanisms of disease. StrVCTVRE runs rapidly and is publicly available.

Introduction

Whole-genome sequencing (WGS) can identify causative variants in clinical cases that elude other diagnostic methods.¹ As the price of WGS falls and it is used more frequently, researchers and clinicians will increasingly observe structural variants (SVs) of unknown significance. SVs are a heterogeneous class of genomic variants that include copy-number variants such as duplications and deletions, rearrangements such as inversions, and mobile element insertions. While a typical short-read WGS study finds 5,000–10,000 SVs per human genome, long-read WGS is able to identify more than 20,000 with much greater reliability.^{2–4} This is two orders of magnitude fewer than the ~3 million single-nucleotide variants (SNVs) identified in a typical WGS study. Still, despite their relatively small number, SVs play a disproportionately large role in genetic disease and are of great interest to clinical geneticists and researchers.^{5,6}

SVs are of clinical interest because they cause many rare diseases. Most SVs identified by WGS are benign, but on average, a given SV is more damaging than an SNV because of its greater size and ability to disrupt multiple exons, create gene fusions, and change gene dosage. In a study of 119 probands who received a molecular diagnosis from short-read WGS, 13% of cases were caused by an

SV.⁷ Similarly, an earlier study that found 7% of congenital scoliosis cases are caused by compound heterozygotes comprised of at least one deletion.⁸ Yet, because SVs continue to be challenging to identify and analyze, these figures may underestimate the true causal role that SVs play in rare disease. Indeed, in some rare diseases, the majority of cases are caused by SVs. For example, deletions cause most known cases of Smith-Magenis syndrome, and duplications cause most known cases of Charcot-Marie-Tooth disease type 1A.⁹ This suggests that for rare disorders, SVs constitute a minor yet appreciable fraction of pathogenic variants.

To continue discovering SVs which cause disease, researchers face a daunting challenge: prioritizing and analyzing the tens of thousands of SVs found by WGS. Best practices for SV prioritization are evolving, and generally mirror steps used to prioritize SNVs. Few SV-tailored impact predictors have been developed, but a small number of published studies have focused on identifying pathogenic SVs from WES^{10,11} and WGS^{7,12,13} and have identified a handful of important steps. Removing low-quality SV calls is essential, as short-read SV callers rarely achieve precision above 80% for deletions and 50% for duplications, even at low recall.¹⁴ Most studies remove SVs seen at high frequency in population databases or internal controls.^{6,15} Moreover, many studies only investigate SVs

¹Biophysics Graduate Group, University of California, Berkeley, Berkeley, CA 94720, USA; ²Center for Computational Biology, University of California, Berkeley, Berkeley, CA 94720, USA; ³Department of Plant and Microbial Biology, University of California, Berkeley, Berkeley, CA 94720, USA; ⁴Department of Biomedical Informatics, Harvard Medical School, Boston, MA 02115, USA; ⁵Division of Genetics, Brigham and Women's Hospital, Boston, MA 02115, USA
*Correspondence: brenner@compbio.berkeley.edu (S.E.B.), sharo@compbio.berkeley.edu (A.G.S.)
<https://doi.org/10.1016/j.ajhg.2021.12.007>

© 2021 The Author(s). This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).



that overlap an exonic region, as non-coding SVs remain particularly difficult to interpret. Depending on its sensitivity, a pathogenic SV discovery pipeline may produce tens to hundreds of rare exon-altering SVs per proband to be investigated. These values are consistent with a recent population-level study that estimates SVs comprise at least 25% of all rare predicted loss-of-function events per genome.¹⁶ Prioritizing SVs will be necessary for the majority of probands, as shown by a study of nearly 500 unresolved cases that found one or more SVs that warranted further investigation in 60% of cases.⁷ Clinically validating all SVs of uncertain significance in a genome is currently infeasible, and cohort size for rare diseases will likely never reach a scale sufficient to statistically associate these SVs with disease. Therefore, computational tools are needed to prioritize and predict the pathogenicity of rare SVs.

Among methods that consider SVs, several annotate the features of SVs but very few prioritize SVs by pathogenicity. General-purpose annotation frameworks such as Ensembl's Variant Effect Predictor (VEP)¹⁷ and SnpEff¹⁸ both annotate SVs with broad consequences on the basis of sequence ontology terms (e.g., transcript_ablation), which we found are not sufficient for effective prioritization. One standalone annotator, SURVIVOR_ant,¹⁹ annotates SVs with genes, repetitive regions, SVs from population databases, and user defined features. This and similar tools put the onus on researchers to provide informative features and determine how to consider these features in combination, a difficult challenge. A complementary approach is to annotate SVs with cataloged SVs known to be pathogenic or benign. One such SV annotator, AnnotSV,²⁰ ranks SVs into five classes on the basis of their overlap with known pathogenic or benign SVs and genes known to be associated with disease or predicted to be intolerant to variation. This approach can be successful when a disease-causing SV has previously been seen in another proband and was cataloged as pathogenic, but we show it has limitations when a disease-causing SV is novel. In contrast, SNVs can be effectively prioritized by methods such as REVEL²¹ and VEST²² that integrate diverse annotations to provide a quantitative score. Similarly powerful methods are needed to predict SV pathogenicity.

In order to provide a summary pathogenicity score to prioritize rare SVs genome-wide, a predictor must address two questions. The first question is whether a gene is most likely associated with a Mendelian phenotype. This relationship can be predicted through gene importance features. The second question is whether an SV impacts gene function, which requires considering intragenic features. Although these are two separate questions, for convenience researchers often combine them into a single summary score. Few methods provide such a summary score for SV pathogenicity. One standalone impact predictor, SVScore,²³ calculates the deleteriousness of all possible SNVs within each SV (via CADD²⁴ scores by default), while considering SV type and gene truncation. SVScore then generates a summary score by aggregating across these

CADD scores (mean of the top 10% by default), and this approach has shown promise in identifying SVs under purifying selection.²³ Another stand-alone predictor, SVFX,²⁵ integrates multiple features into a summary score but focuses on somatic SVs in cancer and germline SVs in common diseases, so we do not discuss it further.

In this manuscript, we introduce StrVCTVRE (structural variant classifier trained on variants rare and exonic), a method that generates a summary pathogenicity score for exon-altering deletions and duplications. We anticipate clinicians and researchers will use StrVCTVRE to prioritize rare SVs associated with Mendelian phenotypes. Since nearly all pathogenic SVs are rare (minor allele frequency [MAF] < 1%), the salient challenge in resolving undiagnosed cases is to distinguish rare pathogenic SVs from rare benign SVs.¹⁶ Existing SV predictors have been trained and assessed on common benign SVs,^{23,25} so they may rely on features that instead separate common SVs from rare SVs and may not be optimal for this clinical question.²⁶ Our unique approach is to train StrVCTVRE to distinguish benign rare SVs from pathogenic rare SVs.

Material and methods

Training, validation, and test datasets

StrVCTVRE was trained on rare SVs from ClinVar,²⁷ gnomAD,¹⁶ and a recent great ape sequencing study.²⁸ StrVCTVRE's performance was evaluated with rare SVs from DECIPHER²⁹ and the 1000 Genomes Project phase 3³⁰ (1KGP).

We retrieved ClinVar SVs²⁷ on January 21, 2020. SVs were retained if they fulfilled all the following requirements: clinical significance of pathogenic, likely pathogenic, pathogenic/likely pathogenic, benign, likely benign, or benign/likely benign; not somatic in origin; type of copy-number loss, copy-number gain, deletion, or duplication; >49 bp in size; at least 1 bp overlap with an exon. We retrieved great ape SVs²⁸ mapped to GRCh38 on April 8, 2019. Deletions were retained if they were absent in humans and homozygous in exactly one of the following species: chimpanzee, gorilla, or orangutan. Only exon-altering deletions > 49 bp were retained. These deletions are subsequently referred to as *apes*. We retrieved gnomAD 2.1.1 SVs¹⁶ (build GRCh37) on June 28, 2019. Only duplications and deletions were retained that were exon altering, >49 bp, and PASS filter. gnomAD SVs were divided into three categories: SVs with a global minor allele frequency (MAF) > 1% ("gnomAD common"), SVs with a global MAF < 1% with at least one individual homozygous for the minor allele ("gnomAD rare benign"), and SVs with a global MAF < 1% with no individuals homozygous for the minor allele ("gnomAD rare unlabeled"). We retrieved GRCh38 "DGV Variants" from the Database of Genomic Variants³¹ release 2016-05-15 on April 08, 2019. MAF of each deletion was calculated as "observedlosses"/(2 × "samplesize"). MAF of each duplication was calculated as "observedgains"/(2 × "samplesize"). Only exon-altering SVs > 49 bp were retained. Those SVs with an MAF greater than 1% are subsequently referred to as "DGV common." We retrieved DECIPHER CNVs (build GRCh37) on Jan 27, 2020. Only exon-altering SVs > 49 bp with pathogenicity of "pathogenic," "likely pathogenic," "benign," or "likely benign" were retained. We only considered benign or likely benign SVs without "full" or

“partial” contribution to disease phenotype. These benign and likely benign SVs were included in all three of the following sets. Set 1 pathogenic SVs consisted of pathogenic or likely pathogenic SVs with “full” contribution to disease phenotype (referred to as “sufficient” in this manuscript). Set 2 SVs consisted of pathogenic or likely pathogenic SVs with “full” or “partial” contribution. Set 3 SVs consisted of pathogenic or likely pathogenic SVs with “full,” “partial,” or “unknown” contribution. Identical SVs with conflicting pathogenicity were removed. SVs were then sorted by size (ascending) and SVs with a reciprocal overlap > 90% were removed; only the first SV was kept. We retrieved 1KGP merged SVs³⁰ on October 22, 2019. Only exon-altering deletions and duplications with a global allele frequency less than 1% were used in our analysis. For the above training and testing SVs, we retrieved SVs mapped to GRCh38 unless noted otherwise. When only SVs mapped to GRCh37 were available, we converted to GRCh38 by using the University of California, Santa Cruz (UCSC) liftover tool.³²

To identify exon-altering SVs, we used exon boundaries from Ensembl biomart,³³ genes v96, GRCh38.p12, limited to genes with HGNC Symbol ID(s) and APPRIS annotation.³⁴ For each gene, a single principal transcript was used, based on the highest APPRIS annotation. For transcripts that tied for highest APPRIS annotation, the longest transcript was used. Exon overlap was determined with bedtools intersect.

To remove near-duplicate SVs in our training and testing data, we performed extensive deduplication of data as follows. Deletions and duplications were considered separately. We ordered benign SVs ($n = 23,239$) (ClinVar benign, ClinVar likely benign, apes, gnomAD rare benign, gnomAD rare unlabeled) and removed duplicates (reciprocal overlap of 90% or greater), keeping the first appearance of an SV. This removed 577 SVs from ClinVar benign/likely benign, five SVs from apes, and 408 SVs from gnomAD. The retained data are subsequently referred to as “benign.” To deduplicate pathogenic SVs ($n = 8,378$), we considered deletions and duplications separately. Exact matches between ClinVar pathogenic and ClinVar likely pathogenic were removed from likely pathogenic. SVs were then sorted by size (ascending). SVs with >90% reciprocal overlap were removed, and the smallest SV was kept. This removed 2,421 pathogenic SVs. The retained data are subsequently referred to as “pathogenic.” Next, exact matches between the benign and pathogenic datasets were removed from both datasets. Finally, duplicates between pathogenic and benign (reciprocal overlap of 90% or greater) were removed from the pathogenic dataset. This removed three benign SVs and 82 pathogenic SVs.

We processed data as follows to ensure we trained only on rare SVs. Pathogenic and benign SVs that exactly matched a DGV common SV were removed. Pathogenic and benign SVs with reciprocal overlap > 90% with an SV in gnomAD common were removed. This removed 30 benign SVs and one pathogenic SV. SVs between 50 bp and 3 Mb were retained and all others were removed.

We found some evidence of acquisition bias in ClinVar data due to the SV size sensitivity of different SV detection methods (see [results](#)). To ensure StrVCTVRE was not learning on this acquisition bias, we matched the size distribution of benign and pathogenic SVs by using the following procedure. After filtering as described above, we organized benign SVs into five tiers: ClinVar likely benign, ClinVar benign, apes, gnomAD rare benign, and gnomAD rare unlabeled. Each pathogenic SV was then matched by size and type (DEL or DUP) to a benign SV, iterating through

each tier. Specifically, each pathogenic SV of size N seeks a benign SV of the same type in the bin $[N - (N/\alpha + 20), N + (N\alpha + 20)]$ where $\alpha = \sqrt[10]{10^6}$ (this bin size derived from Ganel et al.²³). A pathogenic SV first seeks a benign SV in the first benign tier. If matched, the pathogenic and benign SVs are included in the training set, and the benign SV cannot match any further pathogenic SVs. If no match is found in the first benign tier, the same process is repeated while progressing through further benign tiers. Pathogenic SVs that do not find a match in any benign tier are not included in the final training set. This process was continued for all pathogenic SVs, and the resulting data are shown in [Figure S4](#).

After SVs were annotated with features (see below), we identified groups of SVs with identical features, considering pathogenic and benign SVs separately. We removed all but one of these feature-identical SVs in order to avoid overfitting. This removed 37 SVs from the pathogenic training set and 31 SVs from the benign training set. For feature-identical SVs that were present in both the pathogenic and the benign datasets, all feature-identical SVs were removed. This removed 13 SVs.

Structural variant impact predictors

We retrieved VEP¹⁷ v96 on April 16, 2019 and used it to annotate SVs with transcript consequence sequence ontology terms. We retrieved SVScore²³ v0.6 on June 16, 2019. It was run with CADD²⁴ v1.3, which we retrieved on June 16, 2019 by using default settings. We retrieved AnnotSV²⁰ v2.3.2 on Feb 27, 2020. AnnotSV was run with human annotation and default settings. We retrieved X-CNV³⁵ on September 27, 2021, and it was run with default settings on variants converted to GRCh37 via the UCSC liftover tool. We retrieved CADD-SV³⁶ v1.0 on September 13, 2021, and it was run with default settings.

Structural variant features

All gene and exon boundaries used to determine features came from Ensembl Genes v96 as described above. Each SV was annotated with the 17 features listed in [Table 1](#). Expression features were derived from transcript data downloaded from the GTEx Portal v7.³⁷ Exon expression was calculated for each nucleotide as the sum of the transcripts per million (TPM) of fragments that map to that nucleotide. Exon inclusion estimated the proportion of transcripts generated by a gene that include a given nucleotide and was calculated for each nucleotide as the TPM of fragments that map to that nucleotide divided by the sum of TPM that map to the gene containing that nucleotide. For both features, adjacent base pairs with the same value were merged together into genomic intervals. For SVs that overlapped more than one of these genomic intervals, we calculated exon expression by averaging the 400 highest exon expression genomic intervals contained in that SV. The same was done for exon inclusion. All GTEx tissues were used in this analysis.

To determine which conservation feature to use, we assessed the accuracy of both PhastCons³⁹ and PhyloP⁴⁰ in discriminating between pathogenic and benign SVs by using the average of the highest-scoring 200, 400, 600, 800, and 1,000 nucleotides ([Figure S1](#)). The test set consisted of 200 small (<800 bp) SVs randomly selected from our pathogenic and benign SV training datasets (as described above). We found the mean PhyloP score of the 400 most conserved nucleotides in an SV was among the

Table 1. Features used in StrVCTVRE

Feature category	Feature description	Data type	Aggregation method for multiple genes
CDS	fraction of CDS adjacent to start codon that is not disrupted by SV	float	min
CDS	fraction of CDS adjacent to stop codon that is not disrupted by SV	float	min
CDS	fraction of CDS overlapping SV	float	max
conservation	average phyloP score of the 400 most conserved overlapping nucleotides	float	N/A
expression	exon expression (see material and methods)	float	N/A
expression	exon inclusion (see material and methods)	float	N/A
expression	TAD boundary strength (according to Gong et al. ³⁸)	float	max
gene importance	LOEUF of gene	float	min
gene importance	LOEUF of gene where stop codon overlaps SV or > 50% of CDS overlaps SV	float	min
gene importance	pLI of gene	float	max
gene importance	pLI of gene where start codon overlaps SV or > 50% of CDS overlaps SV	float	max
other	all overlapped exons can be skipped in frame	boolean	N/A
other	any overlapped exon is constitutive	boolean	N/A
other	minimum exon transcript order ^a	integer	min
other	number of exons in canonical transcript of gene	integer	min
other	number of exons SV overlaps by 1 or more bp	integer	max
other	SV is deletion or duplication	boolean	N/A

^aExon transcript order was defined as the number of exons preceding a given exon in a gene.

highest accuracy predictors. For both conservation and expression features, if the total overlap between the SV and all exons was less than 400 intervals, then we averaged together the values of the overlapped intervals to calculate the feature. We used median imputation to fill in missing feature annotations.

In our feature correlation analysis, features were clustered by correlation with the linkage and fcluster functions from the SciPy⁴¹ v1.1.0 hierarchical clustering package. The input to this analysis were the features for all SVs used as training data. We reversed values for some features to ensure most matrix correlations are positive.

Random forest classification

StrVCTVRE was implemented as a random forest classifier in Python with scikit-learn⁴² v0.17 with class RandomForestClassifier. We performed a grid search to find the optimal hyperparameters by using a leave-one-chromosome-out cross validation strategy and validation only on ClinVar data, as described previously. The hyperparameters searched included the max depth of a tree (5, 10, 15, no limit), max features considered at each split (1, 2, 3, 4), the minimum samples at each leaf node (1, 2, 4), the minimum samples required to split a node (2, 4), the number of trees generated (500, 1,000, 3,000), and whether to use out-of-bag samples to estimate accuracy (true, false). Several combinations of features performed similarly well, and we chose one that performed well while unlikely to over-fit to the training data—max depth, 10; max features considered at each split, 1; minimum samples at each leaf node, 2; minimum samples required to split a node, 4; number of trees, 1,000; out of bag samples, false. Feature impor-

tance used in figures is also known as Gini importance⁴³ and was calculated with the feature_importances_ attribute of RandomForestClassifier.

Statistical methods used in figures

In Figure 1B, we derived 95% confidence intervals (CIs) by generating 1,000 random forest predictors. In Figure 2A, we generated the data by using a leave-one-chromosome out approach that included all chromosomes besides chromosomes 1, 3, 5, and 7 (e.g., SVs in chromosome 2 were assessed with training data from chromosomes 4, 6, 8, 9, 10, etc.). In Figure 2B, to create the inset testing set, we began with the benign and pathogenic datasets as described above and only retained ClinVar SVs from each dataset. Next, we removed any SVs larger than 3 Mb, and for both the benign and pathogenic dataset, we randomly sampled SVs without replacement, such that SVs were retained if they did not overlap any of the same genes as a previously sampled SV. This resulted in a reduced dataset for both pathogenic and benign SVs, in which every gene was overlapped by at most a single SV. Pathogenic and benign SVs from these reduced datasets were then matched by size as described above, and only results from testing on SVs on chromosomes 1, 3, 5, and 7 are shown in the Figure 2B inset. We note that area under the receiver-operating characteristic curves (AUCs) are often problematic for the evaluation of single-nucleotide or missense impact predictors because of the vast imbalance between pathogenic and benign variants. This imbalance requires methods with adequate specificity. Due to the much smaller number of rare SVs in a genome, receiver-operating characteristic (ROC) plots can be used with care, but with attention still to the leftmost

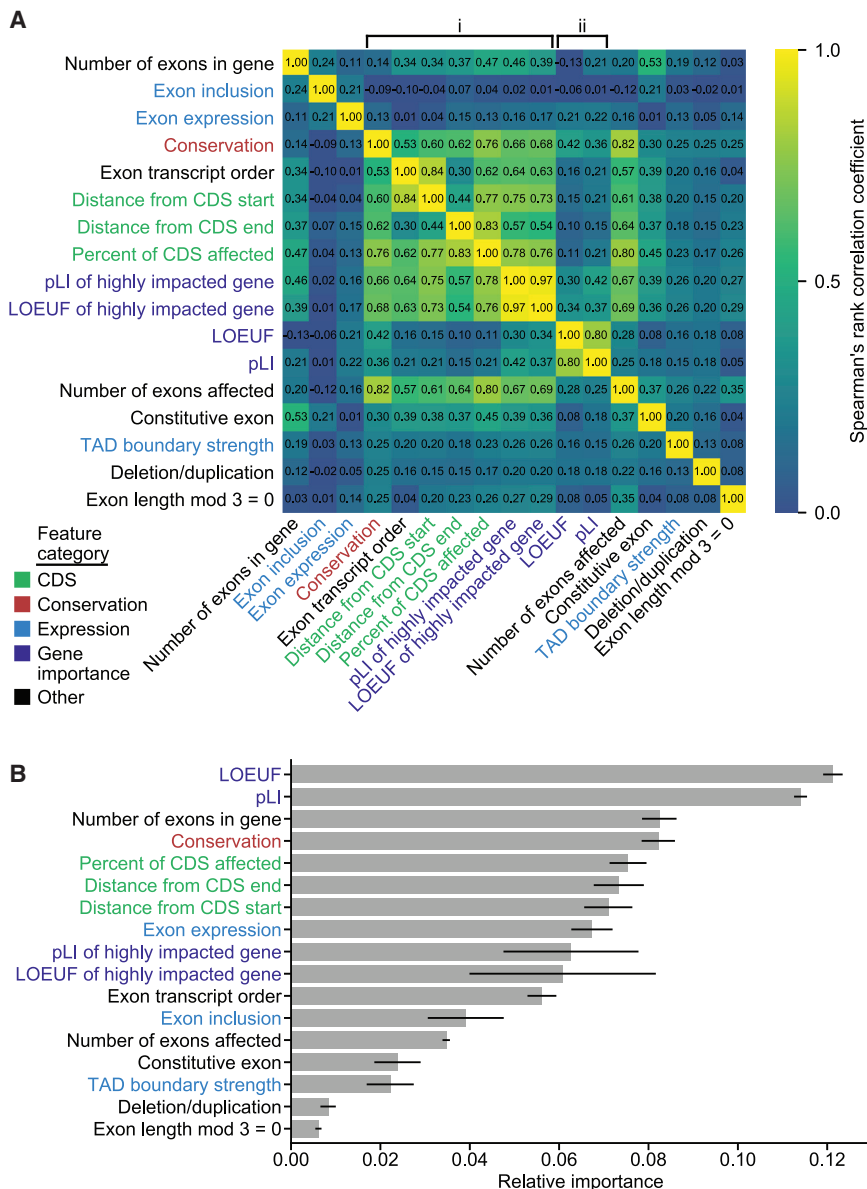


Figure 1. Feature clustering and importance identifies those features providing unique and predictive information

(A) Correlation matrix of StrVCTVRE features in training data. Features were ordered by hierarchical clustering, and some values were reversed to reduce negative correlation between features. Values represent Spearman's rank correlation between features. Text is colored by feature category.

(B) Feature importance of StrVCTVRE features. Gray bars indicate feature importance, estimated with mean decrease in impurity (Gini importance). Black lines indicate 95% confidence intervals. Note that exon expression had high importance yet was uncorrelated with all other features, suggesting it captures unique and predictive information.

dataset of known pathogenic and benign SVs, and the decision nodes are optimized for accuracy. To promote diverse trees, each node of the decision tree uses only a random subset of the features. Finally, StrVCTVRE is assessed on a held-out test dataset and independent test datasets.

Characterization of StrVCTVRE features

To classify SVs, StrVCTVRE employs 17 features in five categories: gene importance, conservation, coding sequence, expression, and exon structure of the disrupted region (see [material and methods](#), [Table 1](#) for details). We assessed gene importance by using two features that summarize the degree of depletion of predicted loss-of-function (pLoF) variants in healthy individuals: pLI⁴⁵ and LOEUF.¹⁵ Although LOEUF is effectively an updated, continuous version of pLI, and the two are highly correlated, we found better performance when both were included rather than just one. To explicitly capture when an important gene is highly impacted by an SV, we included two additional features: pLI of a highly impacted gene and LOEUF of a highly impacted gene. We define a gene as highly impacted when an SV overlaps the APPRIS³⁴ principal start codon or 50% of the coding sequence (CDS). To specifically model CDS disruptions, we used three coding features: percentage of the CDS overlapped by the SV, distance from the CDS start to the nearest position in the SV, and distance from the CDS end to the nearest position in the SV. We included a single conservation feature, phyloP of 100 vertebrates,⁴⁰ by considering the average of the 400 most conserved sites in the SV. PhyloP produced the best classification among

regions reflecting specificity, rather than the AUC computed over the entire curve. We derived AUC 95% confidence intervals by calculating the AUC standard error following Hanley and McNeil.⁴⁴ In the results section “StrVCTVRE performance is higher when assessed on more reliably classified data,” 90% sensitivity thresholds were derived from StrVCTVRE and SVScore performance on the ClinVar held-out test set (dotted line, [Figure 2B](#)).

Results

StrVCTVRE design and assessment

StrVCTVRE is implemented as a random forest, in which many decision trees “vote” for whether a given SV is pathogenic. The StrVCTVRE score reflects the fraction of decision trees that “voted” that the SV is pathogenic. The decision trees are shaped by a learning algorithm in which each tree sees thousands of example SVs from a training

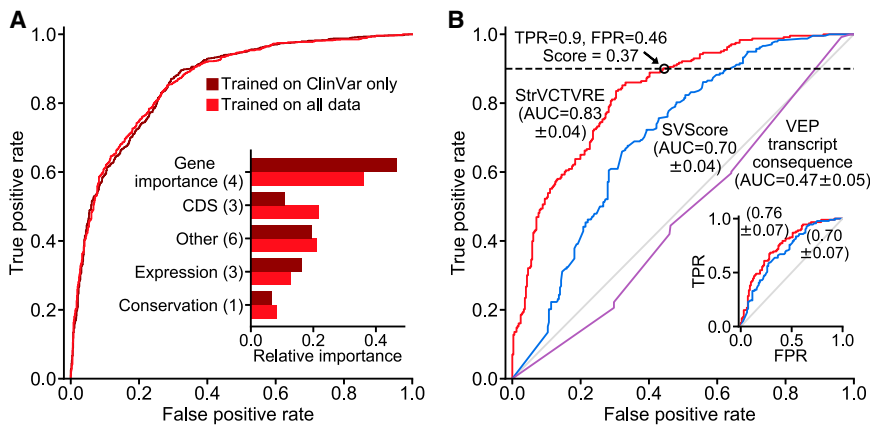


Figure 2. Evaluation of StrVCTVRE on a held-out ClinVar test set and comparison of learned feature importances between training datasets

(A) Receiver-operating characteristic (ROC) comparing StrVCTVRE models trained on two different benign datasets: ClinVar in dark red and all data (ClinVar, SVs common to apes but not humans, and rare gnomAD SVs) in medium red. When tested only on ClinVar data, performance does not significantly differ between the two training sets. However, the feature importances (inset) of the classifier trained on all data (medium red) were more evenly distributed among feature categories. This suggests that unlabeled rare SVs and common ape SVs are a suitable benign training set.

(B) ROC comparing StrVCTVRE (red) to other methods on a held-out test set comprised of ClinVar SVs on chromosomes 1, 3, 5, and 7. Black circle indicates a StrVCTVRE score of 0.37, which we refer to as the ClinVar 90% sensitivity threshold. Inset shows performance on the same held-out test, modified so that each gene is overlapped by a maximum of one SV. AUC with 95% confidence interval is in parentheses.

the conservation features we investigated (see material and methods) and was the most informative conservation feature in a rare missense variant classifier.²¹ To infer expression impacts from the SV, we included the average expression across all tissues for each exon in the SV, the proportion of gene transcripts that included each exon in the SV, and the overlap with known topologically associating domain (TAD) boundaries. To model potential differences that drive the pathogenicity of deletions and duplications, we included as a feature whether an SV is a deletion or duplication. The remaining features were related to the structure of exons in the SV including the number of exons in a disrupted gene, the number of exons disrupted, whether any affected exons were constitutive, whether all disrupted exons could be skipped in frame, and the order of the exon in the transcript. When multiple exons or genes were disrupted, we typically took the value of the most severely impacted one, as appropriate (see material and methods). Missing or non-applicable feature data were replaced by the median value of each feature.

Correlation and relative importance of SV features in StrVCTVRE

Clusters emerged when we calculated these features for our SV training set, computed the correlation between each feature, and clustered by correlation (Figure 1A). The most prominent cluster (labeled i) contains gene importance, conservation, CDS, and one exonic feature. Most correlations were above Spearman's $r = 0.6$. Because both gene importance of highly impacted gene features are present in this cluster, the other features in this cluster may also capture when an important gene is highly disrupted. A smaller cluster (labeled ii) included the remaining gene importance features, pLI and LOEUF. Expression features and deletion/duplication status were the features least correlated with all other features (all $r \leq 0.26$). This low correlation suggests that these features capture unique infor-

mation, which is unsurprising for deletion/duplication status. But given the relative importance of some expression features (Figure 1B), our results suggest expression data contains both orthogonal and valuable information for determining SV pathogenicity. The two features capturing gene importance of a highly impacted gene were the features most correlated with each other ($r = 0.97$), indicating that pLI and LOEUF are generally interchangeable for assessing the importance of highly disrupted genes.

By training on thousands of example SVs, StrVCTVRE discovers which features are useful for discriminating between pathogenic and benign SVs (Figure 1B). Using Gini importance (see material and methods), we found gene importance features were most useful to StrVCTVRE. This was followed by a group of features with similar importance that include the number of exons in a gene, conservation, CDS features, exon expression, and gene importance of a highly impacted gene. The value of these features is largely intuitive; gene importance, CDS, and conservation features are expected to be helpful to assess pathogenicity. In contrast, we suspect number of exons in gene is highly ranked due to sampling bias. We found that many well-studied pathogenic genes have numerous exons (*DMD*, *NF1*, and *BRCA2*), and these genes have many representative SVs in our dataset even after removing near-duplicates (Figure S2, material and methods). This may lead StrVCTVRE to have improved performance on these known clinically relevant genes, but reduced performance genome-wide (discussed further below). Surprisingly, several exonic features had relatively low importance, which may have been caused by the sparsity of SVs in our dataset that alter just a single exon. The low importance of TAD boundaries is counter to findings from a recent cancer SV impact predictor²⁵ and may reflect StrVCTVRE's focus on SVs that impact exons. Additionally, the low importance of deletion/duplication status suggests that on average, for exon-altering deletions and

duplications, the region altered by an SV is more important than whether there was a gain- or loss-of-genome content.

Characterization of StrVCTVRE training and held-out test sets

A total of 7,263 pathogenic or likely pathogenic deletions and 4,551 pathogenic or likely pathogenic duplications were collected from ClinVar,²⁷ a public database of variants cataloged by academic institutions and clinical laboratories. These deletions and duplications include both whole-gene losses and gains and intragenic losses and gains of one or more exons. We restricted our data to deletions and duplications, as they are the only SV types with more than 500 pathogenic examples in ClinVar. Additionally, deletions and duplications constitute the vast majority (>95%) of rare gene-altering SVs.⁶ A set of primarily benign SVs (described in greater detail below) were collected from ClinVar, gnomAD SVs,¹⁶ and a recent great ape sequencing study.²⁸ Because these ape SVs were mapped to the human genome, they may be biased toward more conserved genomic regions. We retained only rare (MAF < 1% in general population) SVs in order to match the challenge faced by SV discovery pipelines. Indeed, 92% of SVs identified through cohort sequencing are rare,¹⁶ so the salient challenge is to distinguish rare pathogenic SVs from rare benign SVs. Existing SV predictors have been trained and assessed on common benign SVs,^{23,25} which may cause them to instead rely on features that separate common from rare SVs and result in lower accuracy in clinical use.²⁶

By training on rare SVs, we intend to achieve better accuracy in the challenge faced in pathogenic SV discovery. To create a rare benign dataset that matches the size range of our pathogenic dataset, we included SVs observed as homozygous at least once in great apes but rare in humans, which we assume should be mostly benign in humans due to our recent shared ancestry with great apes. Our benign dataset also included unlabeled rare SVs from gnomAD SVs. Although we expect a small fraction of these unlabeled SVs are pathogenic, we made two assumptions that mitigated this issue: (1) pathogenic SVs have been depleted by selection, so the large majority of unlabeled SVs are benign, and (2) the fraction of truly pathogenic SVs in the pathogenic and benign training sets is sufficiently different for StrVCTVRE to learn important distinguishing features. By including these additional data sources, we brought the ratio of pathogenic to benign SVs closer to 1:1 in our training set, even at small sizes. This would have been impossible with ClinVar data alone because of the dearth of small benign SVs in ClinVar.

To assess the appropriateness of including SVs from apes and gnomAD in our benign dataset, we explored how performance and feature importance changed with these data included. One predictor was trained only on ClinVar SVs, and a second predictor was trained on ClinVar SVs, ape SVs, and gnomAD SVs (altogether 3.8× more SVs than ClinVar alone). Using leave-one-chromosome-out cross

validation, we found both training sets performed similarly (Figure 2A), supporting our theory that the selected rare unlabeled gnomAD SVs and great ape SVs are sufficiently depleted in pathogenic SVs to be used as a training set of rare, benign SVs. Additionally, the predictor trained on all data showed a distribution of feature importance that is more evenly distributed among feature categories and possibly more robust. This includes a decrease in usefulness of gene importance features, which are likely to be overrepresented in ClinVar data, and an increase in importance in CDS features, which are an important line of evidence for assessing SV pathogenicity.⁴⁶

Before training, we extensively cleansed all data to remove duplicate records within and between datasets, remove common SVs, and remove SVs larger than 3 Mb (see material and methods). Pathogenic deletions and duplications were found to have a large size bias, most likely because of the sensitivity of detection methods to specific size ranges (Figure S3). To avoid training on this acquisition bias, putatively benign SVs were sampled to match the pathogenic SV size distribution (Figure 3; Figure S4). Specifically, in our training data we included only pairs of pathogenic and benign SVs that were of similar size and the same type (deletion or duplication). Using this matching strategy, we were able to include nearly all pathogenic deletions and duplications below 1 Mb. By incorporating ape and gnomAD SVs, we were able to include pathogenic SVs below 10 kilobases (kb), a range nearly absent in ClinVar benign SVs. In the benign training set, 26% of deletions and 75% of duplications came from ClinVar benign or likely benign SVs.

To accurately assess StrVCTVRE's performance, we used a held-out test set of ClinVar SVs on chromosomes 1, 3, 5, and 7 (~20% of the total ClinVar dataset). Only ClinVar SVs were used for testing because it is the highest-confidence dataset. The training set consisted of SVs from all three data sources on all remaining chromosomes. The training set consisted of 2,463 pathogenic SVs and 2,372 benign SVs, and the test set consisted of 244 pathogenic SVs and 334 benign SVs. The test set is of reduced size because pathogenic and benign SVs in the test set were matched on length. None of the SVs in the test set were used to develop the trained algorithm.

StrVCTVRE eliminates more than half of benign SVs from consideration at 90% sensitivity

In discriminating between pathogenic and putatively benign ClinVar SVs in the test dataset, StrVCTVRE performed substantially better than published methods. Performance was measured with the area under the receiver-operating characteristic curve (AUC). The AUC for StrVCTVRE was 0.83 (95% CI: 0.79–0.87). By comparison, SVScore had an AUC of 0.70 (95% CI: 0.66–0.74). StrVCTVRE improved notably in the classification of large duplications and deletions (>1 Mb), a regime in which SVScore by default classifies all SVs as pathogenic (lower left corner of Figure 2B). We also evaluated the predictive

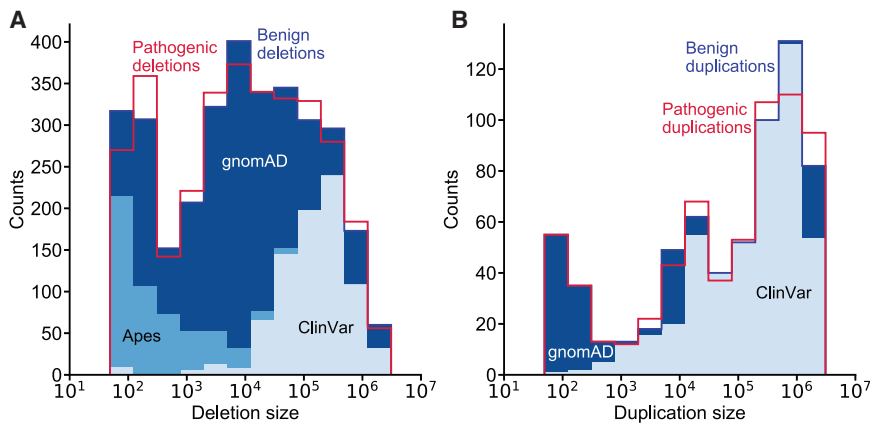


Figure 3. StrVCTVRE draws training data from multiple sources while matching the size distribution of pathogenic SVs

(A and B) Benign training SVs (blue-shaded histograms) closely match the size distribution of pathogenic training SVs (red histogram outlines) and were drawn from multiple datasets. Histogram of pathogenic and benign deletions (A) and duplications (B). (A) Benign deletions are composed of 26% ClinVar, 16% apes, and 58% gnomAD. (B) Benign duplications are composed of 75% ClinVar and 25% gnomAD. We were able to include more small pathogenic SVs in our training data by using apes and gnomAD SVs. Pathogenic SVs are composed entirely of ClinVar pathogenic and likely pathogenic SVs and thus only histogram outlines are shown.

ability of transcript consequence reported by VEP (AUC = 0.47; 95% CI: 0.42–0.52), and we found it performed no better than random. This poor performance was largely due to VEP annotating more benign SVs than pathogenic SVs with its most deleterious sequence ontology term, transcript ablation (Figure S5). The poor performance of transcript consequence from VEP reinforces the known limitations of prioritizing variants with sequence ontology terms in isolation. We also evaluated X-CNV³⁵ and CADD-SV³⁶ on this test dataset (Figure S6). The AUC for X-CNV was 0.68, and the AUC for CADD-SV was 0.70.

As we intend StrVCTVRE to be used to prioritize SVs seen in clinical cases, it needs to perform well in clinically relevant regimes. Clinicians must limit cases in which pathogenic variants are misclassified as benign (false negatives), which requires strong performance at high sensitivity.⁴⁷ When compared to existing methods, StrVCTVRE makes substantial improvements in the high-sensitivity regime, as it is able to capture 90% of pathogenic SVs at a 46% false positive rate (black circle, Figure 2B). StrVCTVRE scores range from 0 to 1, and higher scores indicate a greater likelihood of pathogenicity. In Figure 2B, 90% sensitivity is reached at a StrVCTVRE score of 0.37, which suggests that when used on a collection of SVs called from a clinical cohort, this threshold may identify 90% of pathogenic SVs while reducing the candidate SV list by 54%. We refer to this StrVCTVRE score as the ClinVar 90% sensitivity threshold. StrVCTVRE performed equally well or better on test sets in which duplicates and common SVs were not removed or different size limits were imposed (Figure S7).

We observed apparent clustering in the ClinVar data that led to additional analysis. Genes that are well studied are overlapped by multiple pathogenic SVs cataloged in ClinVar. This resulted in several genes that were over-represented in our test set. Because SVs that overlap the same gene tend to be mostly pathogenic or mostly benign, this results in clustered test data, which may lead to higher variance in AUC performance. While this may yield improved performance for genes of particular interest, it may hide

possible deficits in genome-wide performance. To address this, we randomly generated a test dataset in which each gene is overlapped by at most one SV (Figure 2B inset). We found that the StrVCTVRE AUC was reduced when applied to this dataset, but StrVCTVRE was able to identify pathogenic SVs better than or equal to SVScore at all sensitivities. On this dataset, StrVCTVRE shows a sensitivity of 90% at a false positive rate of 59%. We also considered training StrVCTVRE on a dataset in which each gene is overlapped by at most one SV. We found that this led to a feature importance distribution that is more evenly distributed among feature categories (Figure S8). However, the classifier performance was reduced, so we did not pursue it further.

StrVCTVRE sensitivity threshold is validated on recent clinical SVs

To assess the accuracy of our ClinVar 90% sensitivity threshold and evaluate whether StrVCTVRE performs well on clinical data, we evaluated our method on a set of SVs identified by researchers at the Broad Institute Center for Mendelian Genomics (CMG). These SVs were recently identified through exome sequencing of cohorts with undiagnosed neuromuscular or retinal degeneration disorders.^{48–52} Clinical researchers determined these rare SVs were disease causing or likely disease causing. To avoid overlap between these CMG clinical SVs and StrVCTVRE training SVs, we used a leave-one-chromosome-out approach in which 24 separate StrVCTVRE classifiers were developed, one for each chromosome. For example, CMG clinical SVs on chromosome 1 were predicted by a StrVCTVRE classifier trained on chromosomes 2, 3, 4, etc. The CMG clinical SVs consisted of 32 deletions and two duplications, were located on 14 chromosomes, and had a median size of 12 kb (Table S1). For example, in one proband with a retinal degeneration disorder, a 10 kb deletion on chromosome 19 from 54,121,739 to 54,131,817 received a StrVCTVRE score of 0.85. This variant was confirmed by CMG researchers as causative of the proband's disease.⁵¹ At the ClinVar 90% sensitivity threshold

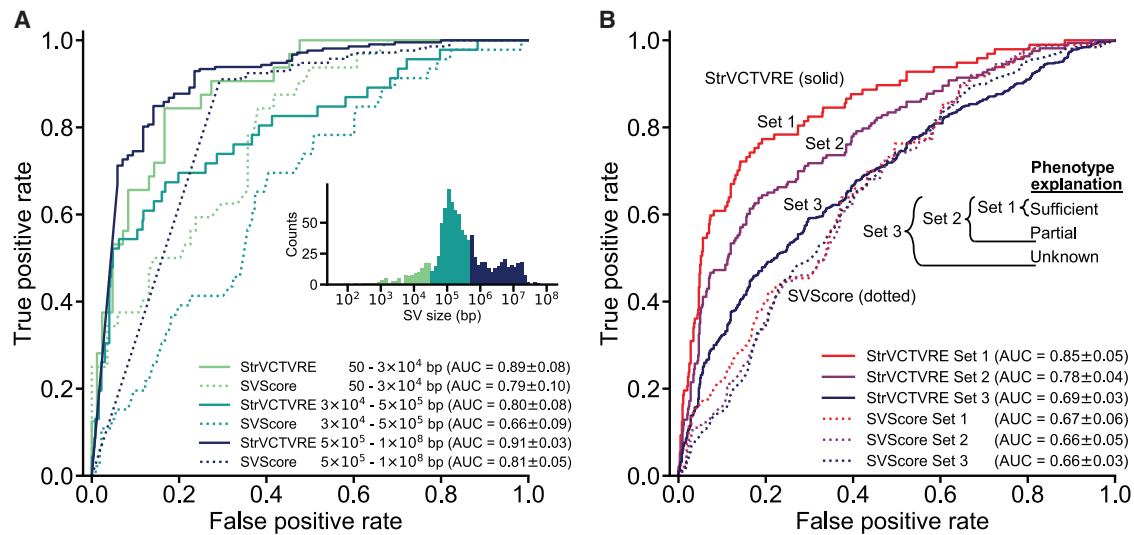


Figure 4. StrVCTVRE performance is consistent across SV size and performance improves on more reliably classified data

(A) Across three size ranges, StrVCTVRE accurately classified variants in an independent test set. In this ROC comparison of StrVCTVRE (solid line) and SVScore (dotted line), three size ranges of SVs were considered. StrVCTVRE performed very well on large and small SVs, while performing well on mid-sized variants. AUC with 95% confidence interval is in parentheses.

(B) When presented with data that are more reliably classified, StrVCTVRE's performance improved. ROC plot showing StrVCTVRE's performance increased as SV contribution to proband phenotype increases from set 3 (includes less confidently classified SVs) to set 2 and from set 2 to set 1 (most confidently classified SVs). The performance of SVScore did not significantly differ between the sets. AUC with 95% confidence interval is in parentheses.

(StrVCTVRE score > 0.37), StrVCTVRE identified 31 of 34 disease-causing SVs (91%) as potentially pathogenic.

Performance of StrVCTVRE on an independent test set from DECIPHER

All held-out test SVs, and a large fraction of training SVs, come from a single database: ClinVar. To independently test StrVCTVRE, we collected pathogenic and benign SVs from DECIPHER, a public database to which clinical scientists submit SVs seen in probands with developmental disorders.²⁹ Because there is some overlap between training ClinVar SVs and DECIPHER SVs, we tested on DECIPHER by using a leave-one-chromosome-out approach, as described above. Additionally, to ensure this DECIPHER test set is independent from our ClinVar test set, we considered only DECIPHER SVs with a reciprocal overlap of less than 10% with any SV used in training or testing StrVCTVRE. This strategy effectively removes any concerns of training and testing on the same or similar SVs. This test set included only DECIPHER variants with the highest classification confidence (set 1, described below). Because StrVCTVRE was trained on SVs less than 3 Mb, and few benign SVs larger than 3 Mb have been observed,³¹ all SVs larger than 3 Mb were scored as pathogenic (given a score of 1). Compared to its performance on the ClinVar test set, StrVCTVRE performed similarly well on the DECIPHER test set, although performance varied across SV size (Figure 4A). On large SVs (>500 kb), StrVCTVRE performed very well (AUC = 0.91; 95% CI: 0.88–0.94; N = 297), partially because most of the SVs larger than 3 Mb are correctly predicted as pathogenic. StrVCTVRE also

performed very well (AUC = 0.89; 95% CI: 0.81–0.97, N = 116) on small SVs (<30 kb), although this is tempered somewhat by the relatively few small SVs in the DECIPHER dataset. StrVCTVRE performed well (AUC = 0.80; 95% CI: 0.72–0.88, N = 545) on mid-length SVs, identifying pathogenic SVs significantly better than SVScore. We also evaluated X-CNV and CADD-SV on these DECIPHER SVs (Figure S9). X-CNV was the second-best performing method but had an AUC below that of StrVCTVRE at all size ranges. CADD-SV, similar to SVScore, performed well on large SVs, but performance was comparably poor on small and mid-length SVs.

StrVCTVRE performance is higher when assessed on more reliably classified data

We expect that some DECIPHER pathogenic SVs are in reality benign. SVs that better explain proband phenotype are more likely to be pathogenic. To investigate the effect of SV pathogenicity on predictor performance, we grouped DECIPHER SVs into three sets. Set 1 consisted of SVs that sufficiently explain the proband phenotype, and these should be reliably pathogenic. Set 2 included SVs that partially explain the proband phenotype and set 1 SVs. Set 3 included SVs with unknown contribution to proband phenotype and set 2 SVs, and therefore, their pathogenicity is less certain. StrVCTVRE was tested with a leave-one-chromosome-out approach. DECIPHER SVs were filtered on the basis of overlap with training and testing data as described above, and only SVs less than 1 Mb in length were retained. We found a consistent trend toward more accurate StrVCTVRE classification in sets

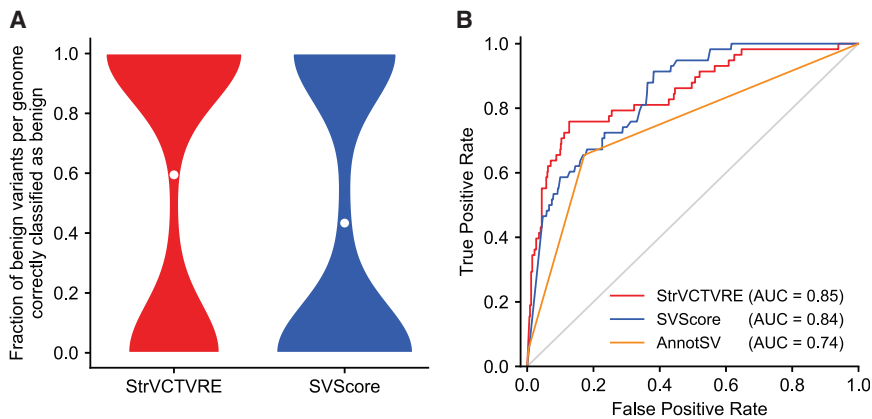


Figure 5. StrVCTVRE helps remove benign SVs from diagnostic consideration and effectively classifies SVs that do not overlap cataloged pathogenic SVs

(A) StrVCTVRE eliminated a significantly larger fraction of benign SVs from consideration than SVScore. When tested on rare exonic SVs from the genomes of 221 putatively healthy individuals, StrVCTVRE was able to correctly classify 59% of putatively benign variants in each genome. White dots represent mean values. For both methods, the threshold for variant consideration was at the ClinVar 90% sensitivity (Figure 2B).

(B) ROC comparing two machine-learning methods with diverse features (StrVCTVRE and SVScore) to one method (AnnotSV) that uses limited features and manually

determined decision boundaries. AnnotSV ranks an SV as “pathogenic” or “likely pathogenic” when the SV overlaps a cataloged pathogenic SV, known disease-associated gene, or gene predicted to be intolerant to variation. To generate this figure, all SVs overlapping any of AnnotSV’s cataloged pathogenic SVs were removed from the DECIPHER set 3 dataset, and the remaining SVs were used for testing. AnnotSV performs relatively poorly on these novel variants. In contrast, the machine-learning methods perform better, possibly because they use more diverse features and have decision boundaries trained on real data. StrVCTVRE scores were generated with a leave-one-chromosome-out approach.

that were more enriched for pathogenic SVs (Figure 4B). This trend was also observed for X-CNV, although to a lesser degree (Figure S10A). However, the same trend was not observed for SVScore nor CADD-SV (Figure S10B). Because StrVCTVRE’s performance improves on presumably more reliably classified data, we have reason to believe StrVCTVRE is making meaningful classifications.

StrVCTVRE eliminates the most benign SVs seen in 221 individuals

Typically, probands with a rare disorder caused by homozygous SVs have one or two pathogenic SVs in their genome, and the remaining SVs are benign. An ideal impact predictor would prioritize the pathogenic homozygous SVs and eliminate from consideration as many of the benign SVs as possible. To evaluate StrVCTVRE’s performance in this scenario, we applied it to SVs called in 2,504 genomes identified by the 1000 Genomes Project phase 3³⁰ (1KGP). Because 1KGP should be depleted of individuals with severe rare disorders, we treated each genome as if it came from a proband with a rare disorder whose pathogenic SVs have been removed. 221 of these genomes had one or more homozygous rare exon-altering SVs, almost all of which should be benign. For each genome, we recorded the fraction of putatively benign SVs that were correctly identified as benign by StrVCTVRE and SVScore (Figure 5A). Since many genomes had just one homozygous exon-altering SV, the distribution is bimodal at 0 and 1. We used our leave-one-chromosome-out predictors (e.g., predicting on 1KGP SVs on chromosome 1 and training StrVCTVRE on all other chromosomes) to score each SV. At the ClinVar 90% sensitivity threshold (StrVCTVRE score > 0.37), on average StrVCTVRE identified 59% of the putatively benign SVs in each genome as benign, compared to 43% when SVScore was used at the

same sensitivity (Wilcoxon paired-rank $p = 8.06 \times 10^{-6}$). In a clinical setting, StrVCTVRE may classify more benign SVs as benign than SVScore, allowing clinicians and researchers to eliminate the most benign homozygous SVs from consideration.

StrVCTVRE performance is reliable even on SVs that do not overlap cataloged pathogenic SVs

Since probands with the same disorder often have SVs altering the same genome element, and recurrent pathogenic *de novo* SVs are known to occur,⁵³ one strategy used to prioritize SVs is to annotate them with overlapping SVs of known pathogenicity. AnnotSV is a popular method to identify pathogenic SVs on the basis of their overlap with both cataloged pathogenic SVs in the National Center for Biotechnology Information’s dbVar. Because it considers cataloged SVs, AnnotSV would most likely perform very well for a proband whose disease-causing SV overlaps a cataloged pathogenic dbVar SV (Figure S11). Yet, many probands have disease-causing SVs that are not cataloged. To address these novel SVs, AnnotSV also considers SV overlap with genes associated with disease or predicted to be intolerant to variation, and it uses manually determined decision boundaries to score SVs (e.g., an SV overlapping a gene with pLI > 0.9 is scored as likely pathogenic). To compare the performance of AnnotSV with machine learning SV impact predictors on novel SVs, we created a dataset of set 3 DECIPHER SVs that do not overlap dbVar SVs used by AnnotSV, and we recorded the prediction accuracy of each method (Figure 5B). AnnotSV performed notably worse on these uncatalogued SVs. We tested StrVCTVRE (with the leave-one-chromosome-out approach) and SVScore on these uncatalogued SVs, and both methods showed significant predictive power, which we attribute to their consideration of features beyond gene intolerance (such as conservation and expression features) and

their use of methods that learn decision boundaries based on training data rather than manually determined boundaries.

Interpreting StrVCTVRE scores

StrVCTVRE scores range from 0 to 1, reflecting the proportion of decision trees in the random forest that classify an SV as pathogenic. Note that StrVCTVRE scores are not probabilities. Although we used the ClinVar 90% sensitivity threshold for evaluation, we advise against using StrVCTVRE scores as a threshold. We instead recommend that greater consideration be given to SVs with greater StrVCTVRE scores. However, thresholds are currently required for computational tools when SVs are classified with the guidelines for sequence variant interpretation recommended by the American College of Medical Genetics and Genomics (ACMG; criteria PP3, BP4).^{46,54} Within the ACMG framework, StrVCTVRE can be used as supporting evidence because it uses multiple lines of computational data. We suspect that higher levels of evidence (e.g., moderate) may be achievable, as shown by Tavtigian et al.⁵⁵ However, when using StrVCTVRE at higher levels of evidence, users should be careful not to also count other ACMG criteria that StrVCTVRE already incorporates, which could lead to double counting. Alternatively, to resolve concerns of double counting, StrVCTVRE can be used just to prioritize variants but not used as evidence. Users then can manually classify SVs of interest by using the full ACMG criteria.

Discussion

As genome sequencing becomes more accessible, clinicians and researchers face a challenge in identifying pathogenic SVs in the thousands identified by sequencing. The ACMG recently offered guidelines for classifying SVs, acknowledging that classification is complex and many pathogenic SVs will be classified as variants of uncertain significance as a result of incomplete knowledge.⁴⁶ SV impact predictors can address this challenge, but few SV impact predictors exist. Although SVs comprise a significant fraction of the loss-of-function mutations that cause rare disease, fewer than 10,000 pathogenic SVs have been cataloged in ClinVar. These SVs have distinct biases toward certain genes and lengths, which leads to acquisition bias that hinders predictor development. Additionally, it is not clear which features are most useful when classifying SVs and how to address the large size range of SVs. StrVCTVRE was developed to address these problems by predicting the impact of exon-altering deletions and duplications in rare genetic disorders. We overcame data limitations and bias by combining SVs from multiple data sources as well as matching pathogenic and benign SVs by size. Because clinicians and researchers must recognize SVs that cause disease among dozens of rare exon-altering SVs detected in a proband, we trained only on rare SVs.

Determining whether a single SV is pathogenic requires consideration of numerous features in combination, as demonstrated by the recent ACMG SV guidelines. Independent of these guidelines, our method identified important features in cataloged SVs. Our findings reinforce clinical guidelines, while also highlighting new areas to explore. Both StrVCTVRE and the ACMG guidelines found gene importance and CDS disruptions to be critical for SV interpretation. Additionally, StrVCTVRE highlighted two features not discussed in the guidelines: conservation and expression. We found exon expression in particular is both predictive and poorly correlated with all other features, suggesting it captures distinctive information for determining pathogenicity. More widespread consideration of expression features could be beneficial for SV classification. StrVCTVRE additionally identified features that are not useful to classify exon-altering SVs, such as TAD boundary strength and whether there is a copy gain or loss. This is consistent with the ACMG guidelines, which do not consider TAD boundaries and provide very similar scoring metrics for both copy gain and loss.

Because SVs range from 50 bp to >10 Mb, it is challenging to accurately classify SVs across this range. Benign SVs in ClinVar are mainly >10 kb, but accurate classification of SVs < 10 kb requires training on benign SVs from the same size range. We accomplished this by training on small benign SVs from great apes and gnomAD. When tested on an independent test set, StrVCTVRE performed well at all size ranges. To be helpful in a clinical setting, a method must perform well at moderately high sensitivity. StrVCTVRE satisfies this requirement and was able to remove 57% of homozygous SVs from consideration at a sensitivity of 90% in the 1KGP dataset. This 90% sensitivity threshold was validated with a dataset of recent SVs observed to cause neuromuscular and retinal degeneration disorders. Overall, we found StrVCTVRE outperforms SVScore in most tasks, even though SVScore's underlying approach, CADD, was trained on >1,000-fold more variants. Additionally, whereas StrVCTVRE was often assessed with a leave-one-chromosome-out approach, SVScore could not be readily modified and thus had the benefit of possibly training on data that overlapped the testing SVs.

StrVCTVRE is accessible as a downloadable command line program (see [data and code availability](#)). Whereas SVScore requires users to download an 80 gigabyte (Gb) CADD file, StrVCTVRE only requires a 9 Gb phyloP file. Because there are an intractably large number of possible SVs, each SV must be scored anew (unlike SNVs for which scores can be pre-computed), and this requires efficient scoring methods. StrVCTVRE runs rapidly and annotates 100,000 gnomAD SVs in 3 min, while SVScore annotates the same SVs in 24 h (Figure S12). StrVCTVRE uses 1 Gb of RAM to annotate 20,000 gnomAD variants. RAM usage will vary on the basis of the fraction of SVs that are exonic and their size distribution. We are also working to make StrVCTVRE scores available through dbNSFP⁵⁶ and WGS. Researchers may use our Dryad repository to

retrain StrVCTVRE on updated data (see [data and code availability](#)).

Following existing predictors, StrVCTVRE predicts the pathogenicity of an SV in isolation. Yet human biology complicates this picture through zygosity and dominance. Because zygosity is not reported for most SVs in ClinVar, StrVCTVRE is zygosity naive. Additionally, StrVCTVRE's pathogenic training dataset consists largely of SVs in genes predicted to lead to dominant disorders ([Figure S13](#)). When tested on sets of SVs predicted to lead to dominant or recessive disorders, StrVCTVRE performs similarly on both ([Figure S14](#)). Researchers who suspect a recessive mode of inheritance may need to consider StrVCTVRE scores in tandem with impact predictor scores for SNVs in *trans* in the same gene. Although genes vary in their tolerance of SVs and dominance, we believe a whole-genome approach will be necessary to identify all pathogenic SVs, including those SVs disrupting genes not currently associated with disease. To identify new disease-associated genes, it may be helpful to consider StrVCTVRE scores in tandem with one of the many methods that assess the match between proband phenotype and known/predicted phenotypes for an affected gene.^{58–60}

A method can only be as good as its training data. SV impact predictors are limited by the relatively small number of identified pathogenic and putatively benign SVs, as well as the over-representation of certain genes in the dataset ([Figure S2](#)). While pathogenic ClinVar variants are commonly used to train variant impact predictors, they are known to include misclassified variants.⁶¹ We know of no characterization of the accuracy of SVs in ClinVar, but work investigating pathogenic SNVs suggest at least 90% are pathogenic on the basis of reclassification rates.⁶² 70% of our pathogenic training SVs have at least one review star in ClinVar, indicating they have supporting evidence that further bolsters our confidence in these data. Nonetheless, data limitations almost certainly curtail the ultimate performance of our approach. StrVCTVRE is unable to classify inversions and insertions due to limited data; however, these have been shown to contribute to a minority of the pLoF events caused by SVs.¹⁶ We are hopeful that additional clinical sequencing studies will identify a more diverse range of SVs, which will be cataloged in open resources such as ClinVar and leveraged to develop more accurate models. We look forward to greater non-coding genome annotations, which will expand our understanding and cataloging of pathogenic noncoding SVs, which remain vexing to classify.

Much of the focus in SV algorithms has been on methods to accurately detect SVs. These methods have left clinicians and researchers awash with SVs not previously known. As experimental methods and algorithms advance, SV detection will improve, but SV interpretation will continue to be challenging. StrVCTVRE advances the clinical evaluation of SVs. During genome-sequencing analysis, some cases contain an SV that matches a cataloged pathogenic SV or satisfies the conditions for patho-

genicity set forth in the ACMG SV guidelines. However, these SVs are often not obvious, and StrVCTVRE can be used to quickly bring these SVs to attention. In the many cases in which no SV is immediately promising, StrVCTVRE aids clinicians and researchers in identifying compelling SVs for manual investigation. Then, if a case remains unresolved by manual investigation, SVs highlighted by StrVCTVRE that are in novel disease-associated genes can be directed to experimental exploration. This will empower researchers to identify novel disease-associated genes where haploinsufficiency and triplosensitivity were not previously known causes of disease. Adoption of structural variant impact predictors will enable clinicians and researchers to make the most of these new data to improve both clinical care and our understanding of basic biology.

Data and code availability

All datasets generated and analyzed during this study are available in the Dryad repository: <https://doi.org/10.6078/D1GM63>, with the following exceptions. Data obtained from DECIPHER, for which access was granted for the current study, are not publicly available because of their sensitive nature. Under reasonable request, DECIPHER data can be requested from <https://www.deciphergenomics.org/about/data-sharing>. A subset of the recent CMG clinical SVs found to cause neuromuscular and retinal degeneration disorders have been made publicly available,^{48–52} but the full dataset is not currently publicly available. These data can be made available by Anne O'Donnell on reasonable request. Information for StrVCTVRE is at <https://compbio.berkeley.edu/proj/strvctvre/>. StrVCTVRE scores can be computed by installing StrVCTVRE at <https://github.com/andrewSharo/StrVCTVRE>. The StrVCTVRE RRID is RRID: SCR_021776.

Supplemental information

Supplemental information can be found online at <https://doi.org/10.1016/j.ajhg.2021.12.007>.

Acknowledgments

We thank Anne O'Donnell, Julia Goodrich, Grace Tiao, Katherine Chao, and Isaac Wong for providing rare clinical SVs. We thank Nilah Ioannidis and Peter Sudmant for helpful feedback during the project. We thank Aashish Adhikari for insightful comments in the initial planning of the project. We thank Azza Althagafi for thorough testing of our GitHub resources as well as Lindsey Guan, Reet Mishra, and Ashish Ramesh for early script testing. We thank Véronique Geoffroy and Jean Muller for helpful discussion of an earlier preprint. This study makes use of data generated by the DECIPHER Consortium. A full list of centers who contributed to the generation of the data is available from the DECIPHER website ([web resources](#)) and via email from decipher@sanger.ac.uk. S.E.B. was unable to fully review this manuscript because of injury. A.G.S. was supported by a National Science Foundation Graduate Research Fellowship, grant no. DGE 1752814. This work was also supported by NIH grant P01 AI138962, a research agreement with Tata Consultancy Services, and the Chan Zuckerberg Biohub. Funding for the DECIPHER project was provided by the Wellcome

Trust. Sequencing and analysis were provided by the Broad Institute of MIT and Harvard Center for Mendelian Genomics (Broad CMG) and was funded by the National Human Genome Research Institute, the National Eye Institute, and the National Heart, Lung, and Blood Institute grant UM1 HG008900 and in part by National Human Genome Research Institute grant R01 HG009141. The funders played no role in the study design, analysis, and interpretation of the results nor in writing the manuscript.

Declaration of interests

S.B. also has appointments at the University of California, San Francisco, Lawrence Berkeley National Laboratory, and the Chan-Zuckerberg Biohub. He is on the board of the Human Genome Variation Society, and he is a member of the ClinGen Sequence Variant Interpretation Working Group. S.S. is an associate member at the Broad Institute. He serves as a consultant and a member of the genetics advisory board at NGM Biopharmaceuticals and a consultant at Inari Agriculture.

Received: August 31, 2021

Accepted: December 9, 2021

Published: January 14, 2022

Web resources

AnnotSV, <https://github.com/lgmgeo/AnnotSV>

CADD, <https://cadd.gs.washington.edu/download>

CADD-SV, <https://github.com/kircherlab/CADD-SV>

ClinVar SVs, ftp://ftp.ncbi.nlm.nih.gov/pub/clinvar/tab_delimited/variant_summary.txt.gz

Database of Genomic Variants, <http://dgv.tcag.ca/dgv/app/downloads>

DECIPHER, <https://decipher.sanger.ac.uk/>

gnomAD SVs, <https://gnomad.broadinstitute.org/downloads#v2-structural-variants>

Great ape SVs, ftp://ftp.ebi.ac.uk/pub/databases/dgva/estd235_Kronenberg_et_al_2017/vcf/

1KGP SVs, ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/integrated_sv_map/supporting/GRCh38_positions/

StrVCTVRE, <https://compbio.berkeley.edu/proj/strvctvre/>

SVScore, <https://github.com/lganel/SVScore>

VEP, <https://github.com/Ensembl/ensembl-vep>

X-CNV, <https://github.com/kbvtmd/XCNV>

References

1. Clark, M.M., Stark, Z., Farnaes, L., Tan, T.Y., White, S.M., Dimmock, D., and Kingsmore, S.F. (2018). Meta-analysis of the diagnostic and clinical utility of genome and exome sequencing and chromosomal microarray in children with suspected genetic diseases. *NPJ Genom. Med.* 3, 16.
2. Lappalainen, T., Scott, A.J., Brandt, M., and Hall, I.M. (2019). Genomic analysis in the age of human genome sequencing. *Cell* 177, 70–84.
3. Ebert, P., Audano, P.A., Zhu, Q., Rodriguez-Martin, B., Porubsky, D., Bonder, M.J., Sulovari, A., Ebler, J., Zhou, W., Serra Mari, R., et al. (2021). Haplotype-resolved diverse human genomes and integrated analysis of structural variation. *Science* 372, eabf7117.
4. Beyter, D., Ingimundardottir, H., Oddsson, A., Eggertsson, H.P., Bjornsson, E., Jonsson, H., Atlason, B.A., Kristmundsdottir, S., Mehrlinger, S., Hardarson, M.T., et al. (2021). Long-read sequencing of 3,622 Icelanders provides insight into the role of structural variants in human diseases and other traits. *Nat. Genet.* 53, 779–786.
5. Weischenfeldt, J., Symmons, O., Spitz, F., and Korbel, J.O. (2013). Phenotypic impact of genomic structural variation: insights from and for human disease. *Nat. Rev. Genet.* 14, 125–138.
6. Abel, H.J., Larson, D.E., Regier, A.A., Chiang, C., Das, I., Kanchi, K.L., Layer, R.M., Neale, B.M., Salerno, W.J., Reeves, C., et al. (2020). Mapping and characterization of structural variation in 17,795 human genomes. *Nature* 583, 83–89.
7. Holt, J.M., Birch, C.L., Brown, D.M., Gajapathy, M., Sosonkina, N., Wilk, B., Wilk, M.A., Spillmann, R.C., Stong, N., and Lee, H. (2019). Identification of Pathogenic Structural Variants in Rare Disease Patients through Genome Sequencing. *bioRxiv*. <https://doi.org/10.1101/627661>.
8. Wu, N., Ming, X., Xiao, J., Wu, Z., Chen, X., Shinawi, M., Shen, Y., Yu, G., Liu, J., Xie, H., et al. (2015). TBX6 null variants and a common hypomorphic allele in congenital scoliosis. *N. Engl. J. Med.* 372, 341–350.
9. Stankiewicz, P., and Lupski, J.R. (2010). Structural variation in the human genome and its role in disease. *Annu. Rev. Med.* 61, 437–455.
10. Ascari, G., Rendtorff, N.D., De Bruyne, M., De Zaeytjij, J., Van Lint, M., Bauwens, M., Van Heetvelde, M., Arno, G., Jacob, J., Creyten, D., et al. (2021). Long-Read Sequencing to Unravel Complex Structural Variants of *CEP78* Leading to Cone-Rod Dystrophy and Hearing Loss. *Front. Cell Dev. Biol.* 9, 664317.
11. Zampaglione, E., Kinde, B., Place, E.M., Navarro-Gomez, D., Maher, M., Jamshidi, F., Nassiri, S., Mazzone, J.A., Finn, C., Schlegel, D., et al. (2020). Copy-number variation contributes 9% of pathogenicity in the inherited retinal degenerations. *Genet. Med.* 22, 1079–1087.
12. Wright, C.F., McRae, J.F., Clayton, S., Gallone, G., Aitken, S., FitzGerald, T.W., Jones, P., Prigmore, E., Rajan, D., Lord, J., et al. (2018). Making new genetic diagnoses with old data: iterative reanalysis and reporting from genome-wide data in 1,133 families with developmental disorders. *Genet. Med.* 20, 1216–1223.
13. Sanchis-Juan, A., Stephens, J., French, C.E., Gleadall, N., Mégy, K., Penkett, C., Shamardina, O., Stirrups, K., Delon, I., Dewhurst, E., et al. (2018). Complex structural variants in Mendelian disorders: identification and breakpoint resolution using short- and long-read genome sequencing. *Genome Med.* 10, 95.
14. Kosugi, S., Momozawa, Y., Liu, X., Terao, C., Kubo, M., and Kamatani, Y. (2019). Comprehensive evaluation of structural variation detection algorithms for whole genome sequencing. *Genome Biol.* 20, 117.
15. Karczewski, K.J., Francioli, L.C., Tiao, G., Cummings, B.B., Alfoldi, J., Wang, Q., Collins, R.L., Laricchia, K.M., Ganna, A., Birnbaum, D.P., et al. (2020). The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* 581, 434–443.
16. Collins, R.L., Brand, H., Karczewski, K.J., Zhao, X., Alfoldi, J., Francioli, L.C., Khera, A.V., Lowther, C., Gauthier, L.D., Wang, H., et al. (2020). A structural variation reference for medical and population genetics. *Nature* 581, 444–451.

17. McLaren, W., Gil, L., Hunt, S.E., Riat, H.S., Ritchie, G.R., Thormann, A., Flicek, P., and Cunningham, F. (2016). The ensemble variant effect predictor. *Genome Biol.* *17*, 122.
18. Cingolani, P., Platts, A., Wang, L., Coon, M., Nguyen, T., Wang, L., Land, S.J., Lu, X., and Ruden, D.M. (2012). A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly (Austin)* *6*, 80–92.
19. Sedlazeck, F.J., Dhroso, A., Bodian, D.L., Paschall, J., Hermes, F., and Zook, J.M. (2017). Tools for annotation and comparison of structural variation. *F1000Res.* *6*, 1795.
20. Geoffroy, V., Herenger, Y., Kress, A., Stoetzel, C., Piton, A., Dollfus, H., and Muller, J. (2018). AnnotSV: an integrated tool for structural variations annotation. *Bioinformatics* *34*, 3572–3574.
21. Ioannidis, N.M., Rothstein, J.H., Pejaver, V., Middha, S., McDonnell, S.K., Baheti, S., Musolf, A., Li, Q., Holzinger, E., Karyadi, D., et al. (2016). REVEL: an ensemble method for predicting the pathogenicity of rare missense variants. *Am. J. Hum. Genet.* *99*, 877–885.
22. Carter, H., Douville, C., Stenson, P.D., Cooper, D.N., and Karchin, R. (2013). Identifying Mendelian disease genes with the variant effect scoring tool. *BMC Genomics* *14* (Suppl 3), S3.
23. Ganel, L., Abel, H.J., Hall, I.M.; and FinMetSeq Consortium (2017). SVScore: an impact prediction tool for structural variation. *Bioinformatics* *33*, 1083–1085.
24. Rentzsch, P., Witten, D., Cooper, G.M., Shendure, J., and Kircher, M. (2019). CADD: predicting the deleteriousness of variants throughout the human genome. *Nucleic Acids Res.* *47* (D1), D886–D894.
25. Kumar, S., Harmanci, A., Vytheeswaran, J., and Gerstein, M.B. (2020). SVFX: a machine learning framework to quantify the pathogenicity of structural variants. *Genome Biol.* *21*, 274.
26. Li, M.-X., Kwan, J.S., Bao, S.-Y., Yang, W., Ho, S.-L., Song, Y.-Q., and Sham, P.C. (2013). Predicting mendelian disease-causing non-synonymous single nucleotide variants in exome sequencing studies. *PLoS Genet.* *9*, e1003143.
27. Landrum, M.J., Lee, J.M., Benson, M., Brown, G.R., Chao, C., Chitipiralla, S., Gu, B., Hart, J., Hoffman, D., Jang, W., et al. (2018). ClinVar: improving access to variant interpretations and supporting evidence. *Nucleic Acids Res.* *46* (D1), D1062–D1067.
28. Kronenberg, Z.N., Fiddes, I.T., Gordon, D., Murali, S., Cantsilieris, S., Meyerson, O.S., Underwood, J.G., Nelson, B.J., Chaisson, M.J., and Dougherty, M.L. (2018). High-resolution comparative analysis of great ape genomes. *Science* *360*, eaar6343.
29. Firth, H.V., Richards, S.M., Bevan, A.P., Clayton, S., Corpas, M., Rajan, D., Van Vooren, S., Moreau, Y., Pettett, R.M., and Carter, N.P. (2009). DECIPHER: database of chromosomal imbalance and phenotype in humans using ensembl resources. *Am. J. Hum. Genet.* *84*, 524–533.
30. Sudmant, P.H., Rausch, T., Gardner, E.J., Handsaker, R.E., Abyzov, A., Huddleston, J., Zhang, Y., Ye, K., Jun, G., Fritz, M.H.-Y., et al. (2015). An integrated map of structural variation in 2,504 human genomes. *Nature* *526*, 75–81.
31. MacDonald, J.R., Ziman, R., Yuen, R.K., Feuk, L., and Scherer, S.W. (2014). The Database of Genomic Variants: a curated collection of structural variation in the human genome. *Nucleic Acids Res.* *42*, D986–D992.
32. Kent, W.J., Sugnet, C.W., Furey, T.S., Roskin, K.M., Pringle, T.H., Zahler, A.M., and Haussler, D. (2002). The human genome browser at UCSC. *Genome Res.* *12*, 996–1006.
33. Howe, K.L., Achuthan, P., Allen, J., Allen, J., Alvarez-Jarreta, J., Amode, M.R., Armean, I.M., Azov, A.G., Bennett, R., Bhai, J., et al. (2021). Ensembl 2021. *Nucleic Acids Res.* *49* (D1), D884–D891.
34. Rodriguez, J.M., Rodriguez-Rivas, J., Di Domenico, T., Vázquez, J., Valencia, A., and Tress, M.L. (2018). APPRIS 2017: principal isoforms for multiple gene sets. *Nucleic Acids Res.* *46* (D1), D213–D217.
35. Zhang, L., Shi, J., Ouyang, J., Zhang, R., Tao, Y., Yuan, D., Lv, C., Wang, R., Ning, B., Roberts, R., et al. (2021). X-CNV: genome-wide prediction of the pathogenicity of copy number variations. *Genome Med.* *13*, 132.
36. Kleinert, P., and Kircher, M. (2021). CADD-SV—a framework to score the effects of structural variants in health and disease. *bioRxiv*. <https://doi.org/10.1101/2021.07.10.451798>.
37. Lonsdale, J., Thomas, J., Salvatore, M., Phillips, R., Lo, E., Shad, S., Hasz, R., Walters, G., Garcia, F., Young, N.; and GTEx Consortium (2013). The genotype-tissue expression (GTEx) project. *Nat. Genet.* *45*, 580–585.
38. Gong, Y., Lazaris, C., Sakellaropoulos, T., Lozano, A., Kambadur, P., Ntziachristos, P., Aifantis, I., and Tsirigos, A. (2018). Stratification of TAD boundaries reveals preferential insulation of super-enhancers by strong boundaries. *Nat. Commun.* *9*, 542.
39. Siepel, A., Bejerano, G., Pedersen, J.S., Hinrichs, A.S., Hou, M., Rosenbloom, K., Clawson, H., Spieth, J., Hillier, L.W., Richards, S., et al. (2005). Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res.* *15*, 1034–1050.
40. Pollard, K.S., Hubisz, M.J., Rosenbloom, K.R., and Siepel, A. (2010). Detection of nonneutral substitution rates on mammalian phylogenies. *Genome Res.* *20*, 110–121.
41. Virtanen, P., Gommers, R., Oliphant, T.E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., et al. (2020). SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat. Methods* *17*, 261–272.
42. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., and Dubourg, V. (2011). Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.* *12*, 2825–2830.
43. Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The elements of statistical learning: data mining, inference, and prediction* (Springer Science & Business Media).
44. Hanley, J.A., and McNeil, B.J. (1982). The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology* *143*, 29–36.
45. 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., et al. (2016). Analysis of protein-coding genetic variation in 60,706 humans. *Nature* *536*, 285–291.
46. Riggs, E.R., Andersen, E.F., Cherry, A.M., Kantarci, S., Kearney, H., Patel, A., Raca, G., Ritter, D.I., South, S.T., Thorland, E.C., et al. (2020). Technical standards for the interpretation and reporting of constitutional copy-number variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics (ACMG) and the Clinical Genome Resource (ClinGen). *Genet. Med.* *22*, 245–257.

47. Jagadeesh, K.A., Wenger, A.M., Berger, M.J., Guturu, H., Stenson, P.D., Cooper, D.N., Bernstein, J.A., and Bejerano, G. (2016). M-CAP eliminates a majority of variants of uncertain significance in clinical exomes at high sensitivity. *Nat. Genet.* *48*, 1581–1586.
48. Donkervoort, S., Mohassel, P., Laugwitz, L., Zaki, M.S., Kamsteeg, E.J., Maroofian, R., Chao, K.R., Verschuuren-Bemelmans, C.C., Horber, V., Fock, A.J.M., et al. (2020). Biallelic loss of function variants in SYT2 cause a treatable congenital onset presynaptic myasthenic syndrome. *Am. J. Med. Genet. A.* *182*, 2272–2283.
49. Töpf, A., Johnson, K., Bates, A., Phillips, L., Chao, K.R., England, E.M., Laricchia, K.M., Mullen, T., Valkanas, E., Xu, L., et al. (2020). Sequential targeted exome sequencing of 1001 patients affected by unexplained limb-girdle weakness. *Genet. Med.* *22*, 1478–1488.
50. Ravenscroft, G., Clayton, J.S., Faiz, F., Sivadorai, P., Milnes, D., Cincotta, R., Moon, P., Kamien, B., Edwards, M., Delatycki, M., et al. (2021). Neurogenetic fetal akinesia and arthrogryposis: genetics, expanding genotype-phenotypes and functional genomics. *J. Med. Genet.* *58*, 609–618.
51. Zampaglione, E., Maher, M., Place, E.M., Wagner, N.E., DiTroia, S., Chao, K.R., England, E., Broad, C., Catomeris, A., and Nassiri, S. (2021). The Importance of Automation in Genetic Diagnosis: Lessons from Analyzing an Inherited Retinal Degeneration Cohort with the Mendelian Analysis Toolkit (MATK). *medRxiv*. <https://doi.org/10.1101/2021.04.09.21255188>.
52. Wahlster, L., Verboon, J.M., Ludwig, L.S., Black, S.C., Luo, W., Garg, K., Voit, R.A., Collins, R.L., Garimella, K., Costello, M., et al. (2021). Familial thrombocytopenia due to a complex structural variant resulting in a WAC-ANKRD26 fusion transcript. *J. Exp. Med.* *218*, e20210444.
53. Sanders, S.J., Ercan-Sencicek, A.G., Hus, V., Luo, R., Murtha, M.T., Moreno-De-Luca, D., Chu, S.H., Moreau, M.P., Gupta, A.R., Thomson, S.A., et al. (2011). Multiple recurrent de novo CNVs, including duplications of the 7q11.23 Williams syndrome region, are strongly associated with autism. *Neuron* *70*, 863–885.
54. Richards, S., Aziz, N., Bale, S., Bick, D., Das, S., Gastier-Foster, J., Grody, W.W., Hegde, M., Lyon, E., Spector, E., et al. (2015). Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genet. Med.* *17*, 405–424.
55. Tavtigian, S.V., Greenblatt, M.S., Harrison, S.M., Nussbaum, R.L., Prabhu, S.A., Boucher, K.M., Biesecker, L.G.; and ClinGen Sequence Variant Interpretation Working Group (ClinGen SVI) (2018). Modeling the ACMG/AMP variant classification guidelines as a Bayesian classification framework. *Genet. Med.* *20*, 1054–1060.
56. Liu, X., Li, C., Mou, C., Dong, Y., and Tu, Y. (2020). dbNSFP v4: a comprehensive database of transcript-specific functional predictions and annotations for human nonsynonymous and splice-site SNVs. *Genome Med.* *12*, 103.
57. Liu, X., White, S., Peng, B., Johnson, A.D., Brody, J.A., Li, A.H., Huang, Z., Carroll, A., Wei, P., Gibbs, R., et al. (2016). WGSAs: an annotation pipeline for human genome sequencing studies. *J. Med. Genet.* *53*, 111–112.
58. Köhler, S., Schulz, M.H., Krawitz, P., Bauer, S., Dölken, S., Ott, C.E., Mundlos, C., Horn, D., Mundlos, S., and Robinson, P.N. (2009). Clinical diagnostics in human genetics with semantic similarity searches in ontologies. *Am. J. Hum. Genet.* *85*, 457–464.
59. Singleton, M.V., Guthery, S.L., Voelkerding, K.V., Chen, K., Kennedy, B., Margraf, R.L., Durtschi, J., Eilbeck, K., Reese, M.G., Jorde, L.B., et al. (2014). Phevor combines multiple biomedical ontologies for accurate identification of disease-causing alleles in single individuals and small nuclear families. *Am. J. Hum. Genet.* *94*, 599–610.
60. Zemojtel, T., Köhler, S., Mackenroth, L., Jäger, M., Hecht, J., Krawitz, P., Graul-Neumann, L., Doelken, S., Ehmke, N., Spielmann, M., et al. (2014). Effective diagnosis of genetic disease by computational phenotype analysis of the disease-associated genome. *Sci. Transl. Med.* *6*, 252ra123.
61. Shah, N., Hou, Y.-C.C., Yu, H.-C., Sainger, R., Caskey, C.T., Venter, J.C., and Telenti, A. (2018). Identification of misclassified ClinVar variants via disease population prevalence. *Am. J. Hum. Genet.* *102*, 609–619.
62. Harrison, S.M., and Rehm, H.L. (2019). Is ‘likely pathogenic’ really 90% likely? Reclassification data in ClinVar. *Genome Med.* *11*, 72.

The American Journal of Human Genetics, Volume 109

Supplemental information

**StrVCTVRE: A supervised learning method
to predict the pathogenicity of
human genome structural variants**

Andrew G. Sharo, Zhiqiang Hu, Shamil R. Sunyaev, and Steven E. Brenner

Supplemental Figures

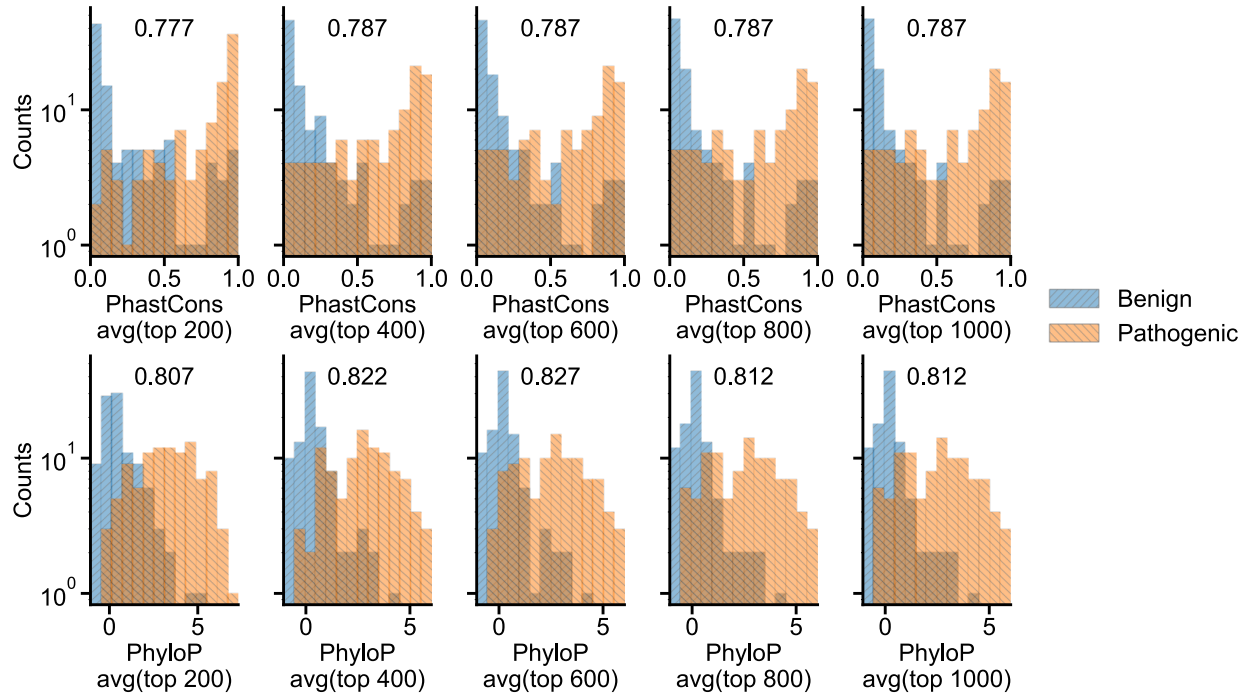


Fig. S1. Average of the top 400 PhyloP scores is one of the best indicators of SV pathogenicity. Comparison of PhyloP and PhastCons in differentiating between pathogenic and benign SVs. The columns represent different values of N, and the top row is PhastCons, while the bottom row is PhyloP. In each plot, the classification accuracy of a linear classifier is shown, and was calculated using a single feature: the average of the largest N PhyloP or PhastCons values among all positions overlapped by the SV. Each plot uses the same underlying SVs and shows a histogram of SV features.

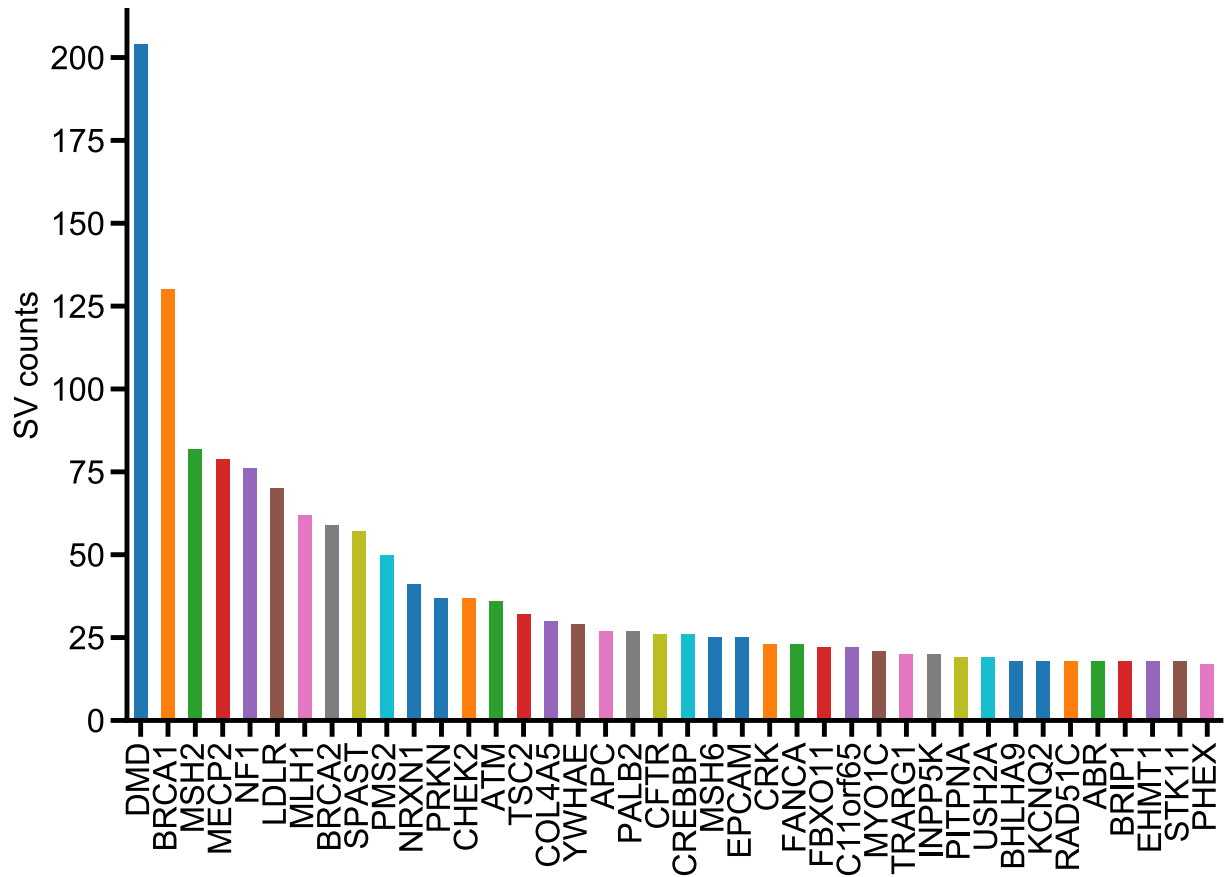


Fig. S2. Certain well-studied genes were overrepresented in our dataset. This plot shows the number of SVs overlapping each gene in the pathogenic ClinVar SVs (only showing the top 40 genes). For context, the median number of SVs overlapping each gene in our training dataset of pathogenic ClinVar SVs is 1.

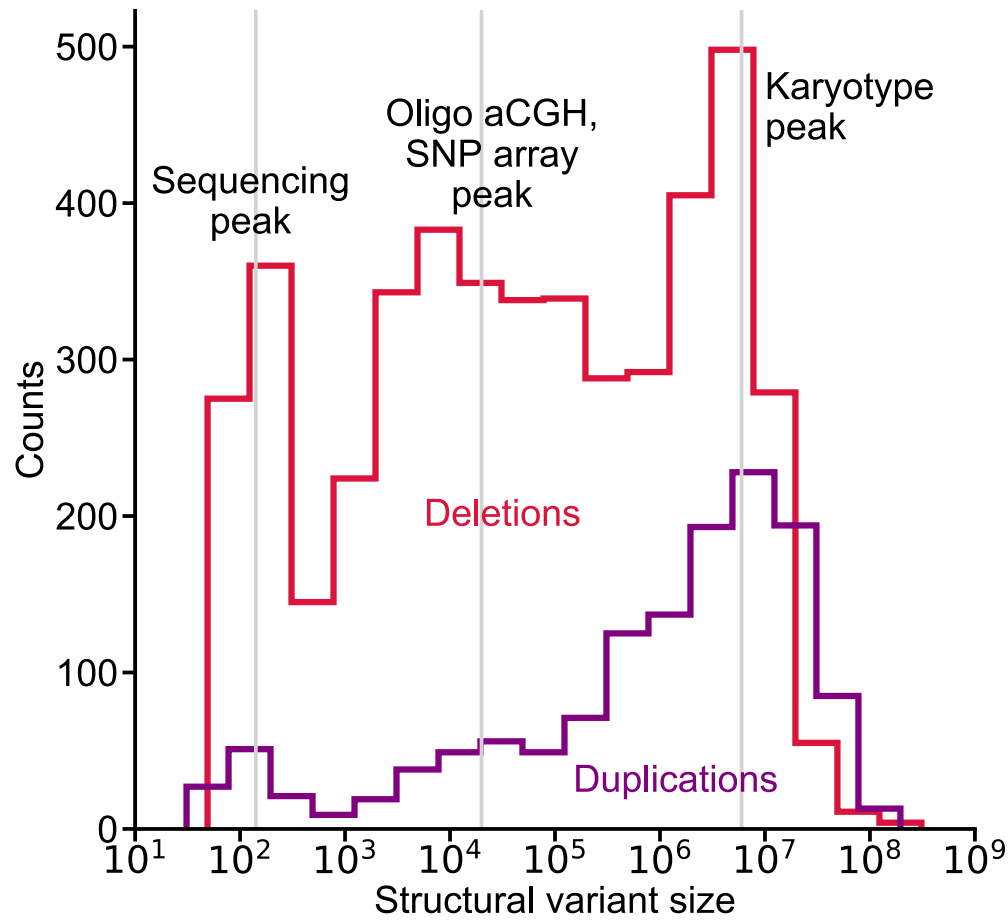


Fig. S3. Histogram of pathogenic SVs shows acquisition bias in SV size. Pathogenic SVs collected from ClinVar display three peaks consistent across both deletions and duplications. These peaks roughly correspond to the sensitivity range for three major methods to identify SVs: sequencing, oligo aCGH/SNP array, and karyotyping.

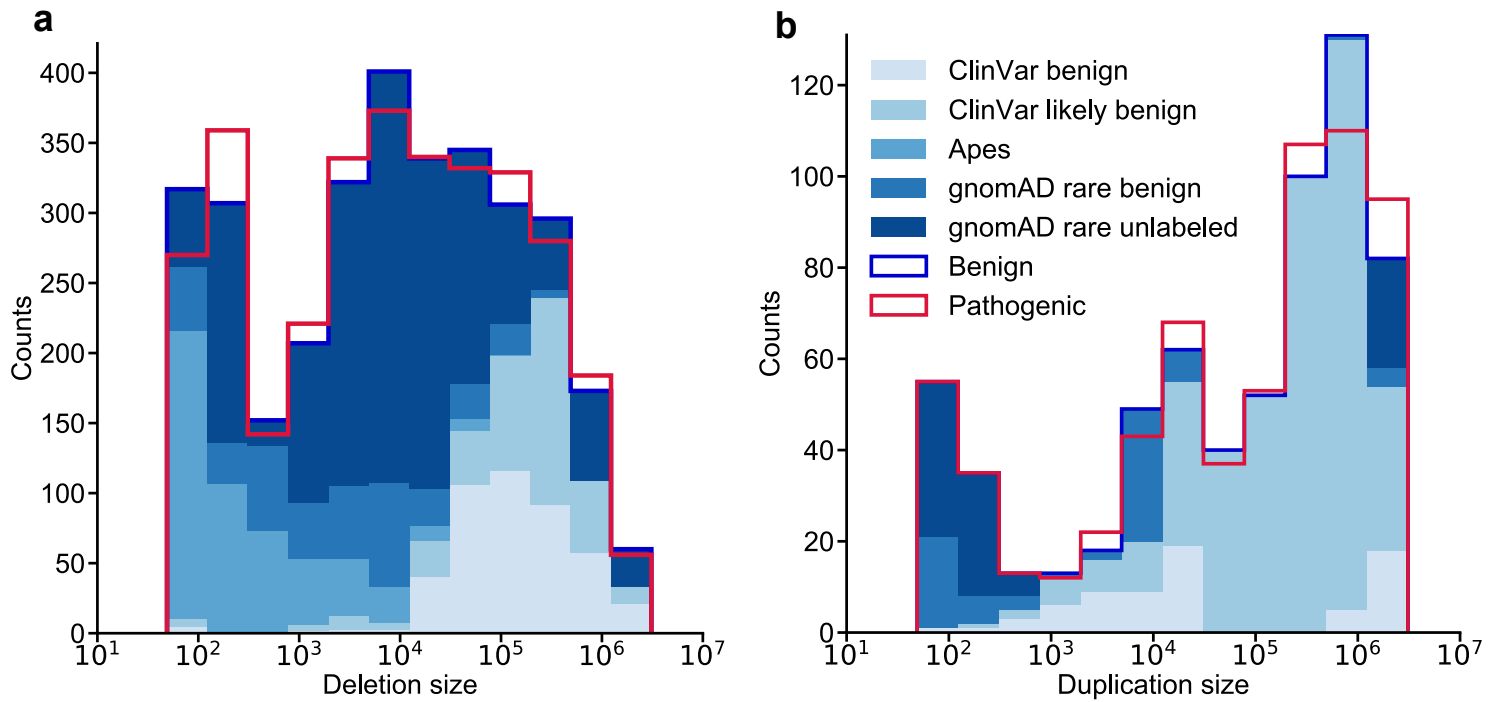


Fig. S4. Drawing SVs from multiple datasets allowed StrVCTVRE to train on a balanced dataset of SVs across a large size range. This is a more detailed version of Fig. 3. **a** Benign deletions were composed of 14% ClinVar benign, 12% ClinVar likely benign, 16% apes, 12% gnomAD rare benign, and 46% gnomAD rare unlabeled. **b** Benign duplications were composed of 11% ClinVar benign, 64% ClinVar likely benign, 11% gnomAD rare benign, and 14% gnomAD rare unlabeled.

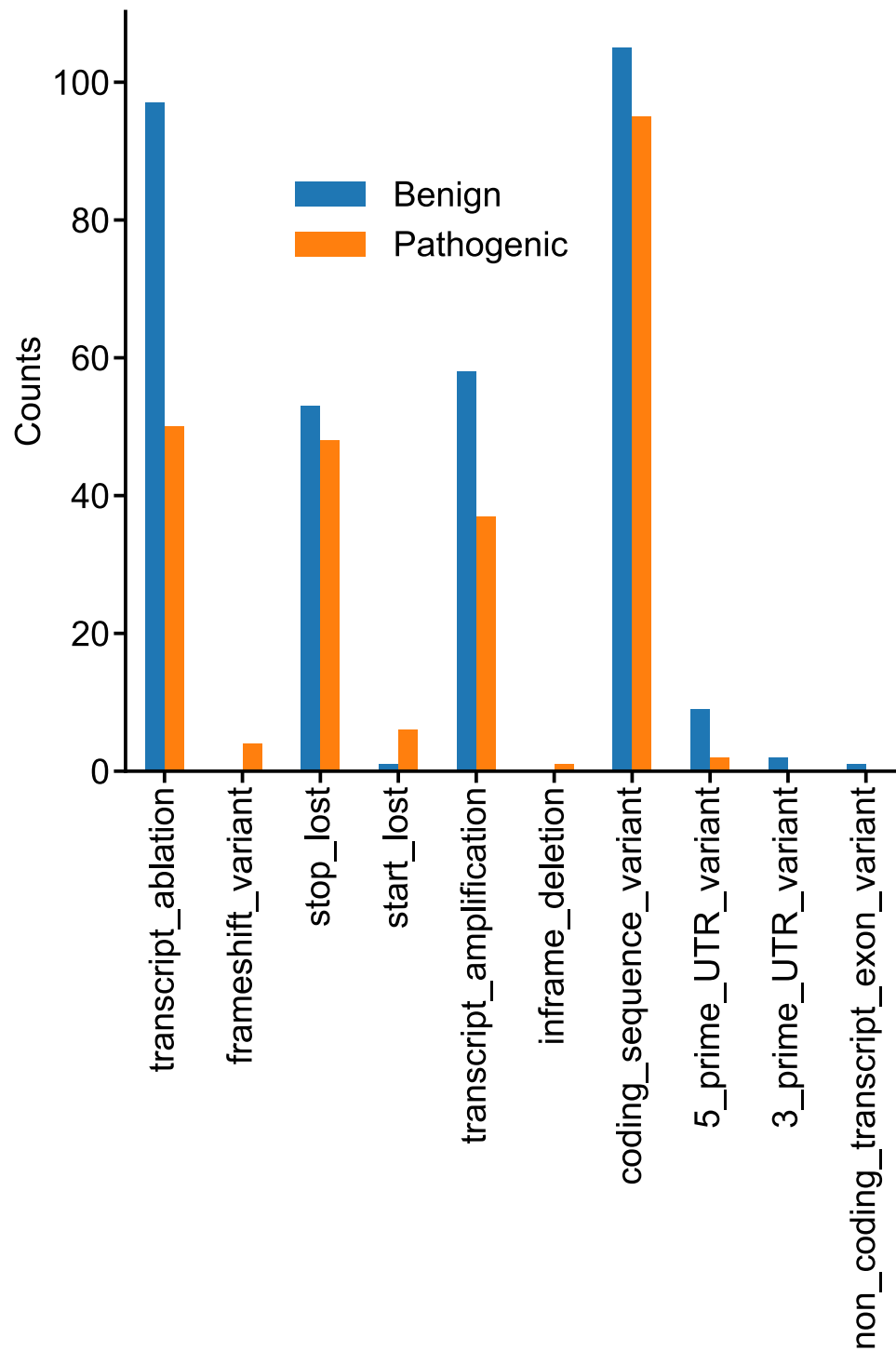


Fig. S5. Number of ClinVar test SVs (chromosomes 1, 3, 5, 7; from Fig. 2b) that VEP annotated with various sequence ontology terms. VEP classified more benign SVs than pathogenic SVs as `transcript_ablation`, leading to its poor performance in SV classification. Categories are ordered from left to right in descending deleteriousness.

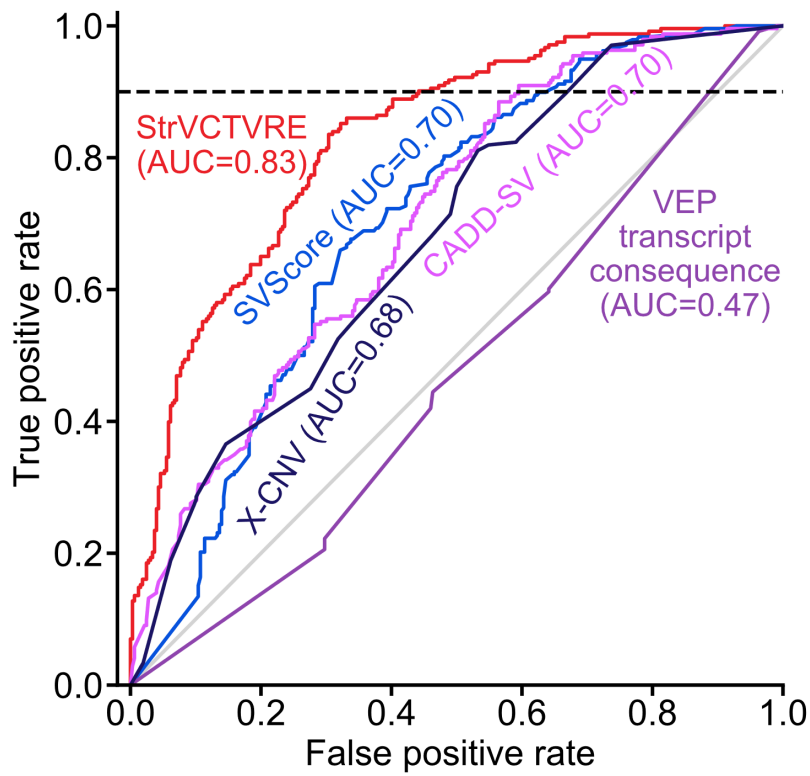


Fig. S6. Extended version of Fig. 2b. Receiver operating characteristic comparing StrVCTVRE (red) to SVScore (medium blue), CADD-SV (pink), X-CNV (dark blue), and VEP (purple) on a held-out test set comprised of ClinVar SVs on chromosomes 1, 3, 5, and 7.

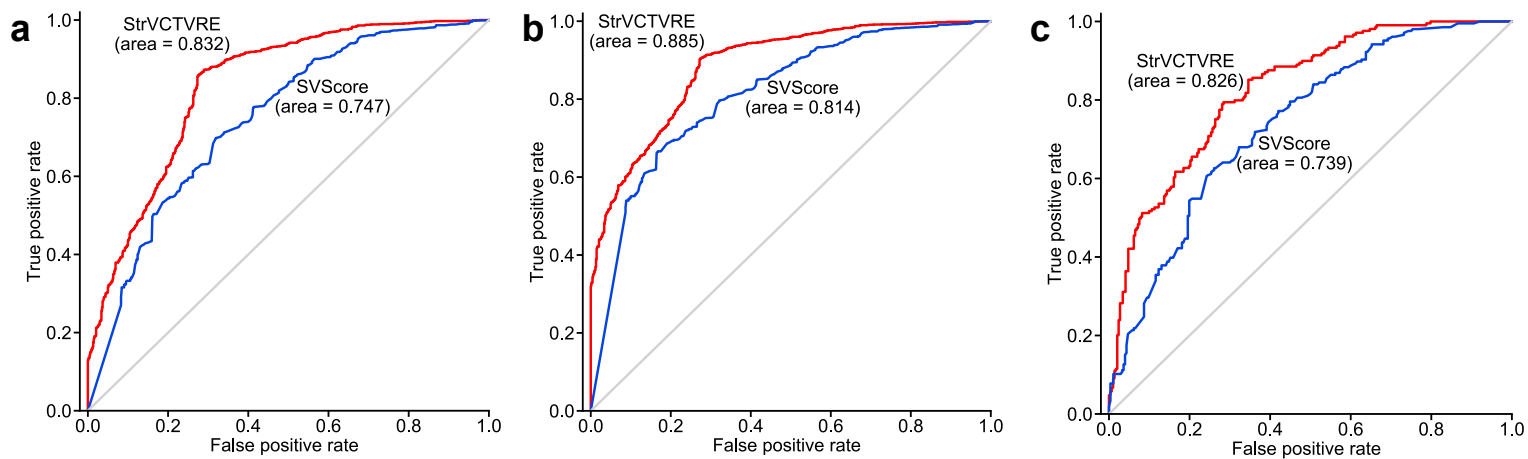


Fig. S7. StrVCTVRE and SVScore performance comparison on various held-out test datasets (ClinVar SVs on chromosomes 1, 3, 5, 7). **a** Duplicates and common variants were not removed from the test data. **b** Duplicates and common variants were not removed from the test data, and no size limit was imposed. **c** Duplicates and common variants were removed from the test data, and only SVs smaller than 1 Mb were included.

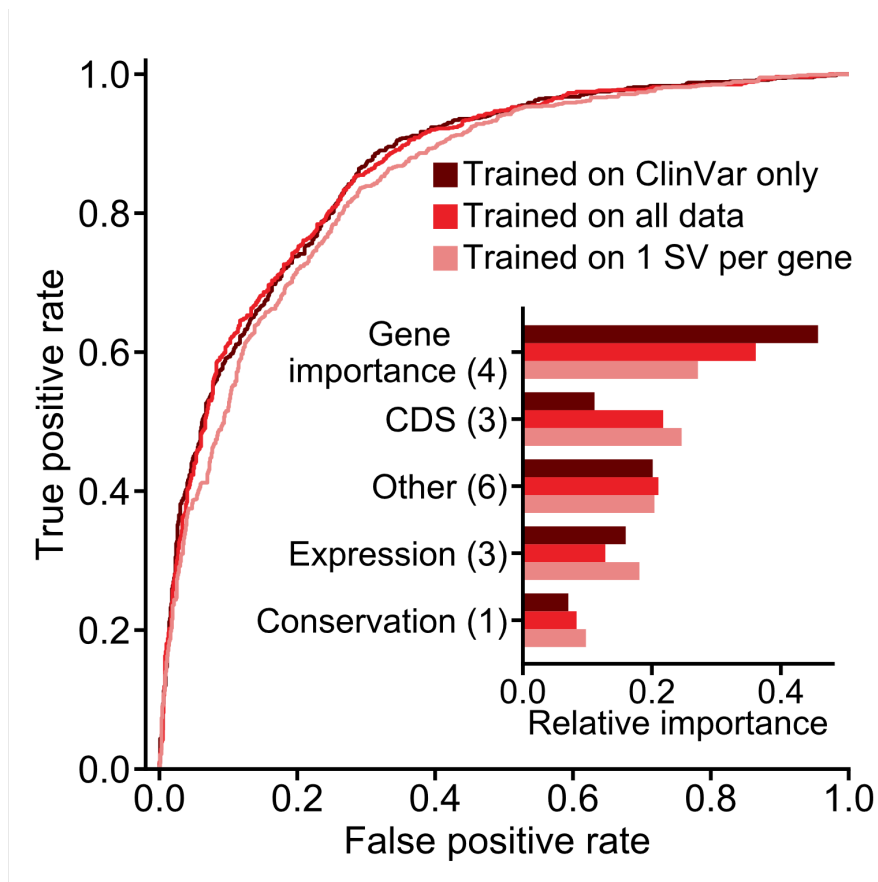


Fig. S8. Receiver operating characteristic comparing StrVCTVRE models trained on three different datasets: ClinVar in dark red, all data (ClinVar, SVs common to apes but not humans, and rare gnomAD SVs) in medium red, and a dataset in which each gene is overlapped by at most one SV in light red. The feature importances of the classifier trained on the one SV per gene dataset (light red) are more evenly distributed than the other two datasets, but the AUC performance is decreased by 0.02 when this classifier is tested on held-out SVs.

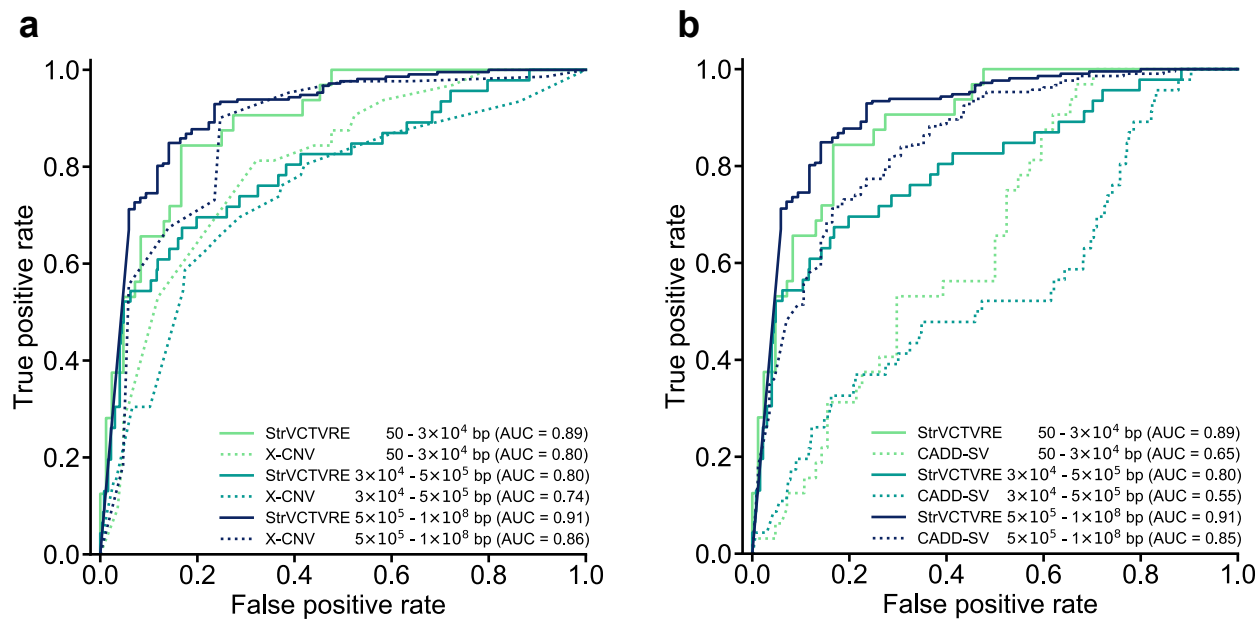


Fig. S9. Extended version of Fig. 4a. Comparison of StrVCTVRE with X-CNV and CADD-SV on a set of DECIPHER SVs across three size ranges. **a** StrVCTVRE (solid lines) performance compared to X-CNV (dotted lines). Line color denotes SV size range. **b** StrVCTVRE (solid lines) performance compared to CADD-SV (dotted lines). Line color denotes SV size range.

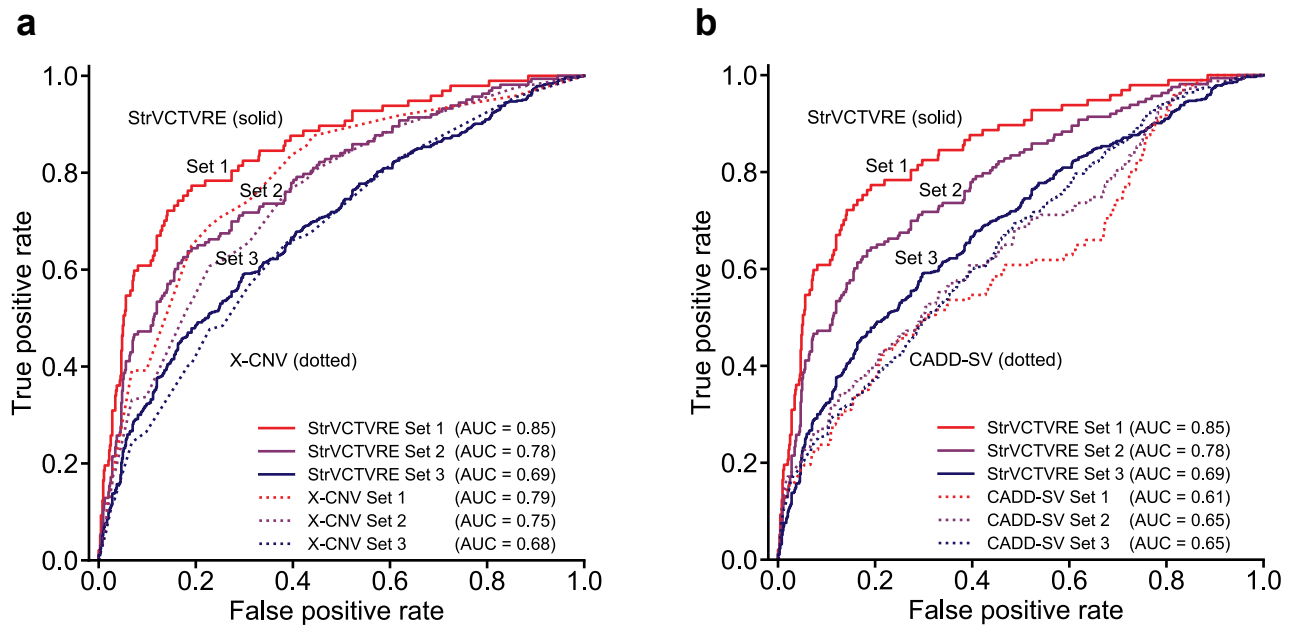


Fig. S10. Extended version of Fig. 4b. Comparison of StrVCTVRE with X-CNV and CADD-SV on a set of DECIPHER SVs with varying levels of contribution to proband phenotype. SV contribution to proband phenotype increases from set 3 (includes less confidently classified SVs) to set 2 and from set 2 to set 1 (most confidently classified SVs). **a** StrVCTVRE (solid lines) performance compared to X-CNV (dotted lines). **b** StrVCTVRE (solid lines) performance compared to CADD-SV (dotted lines).

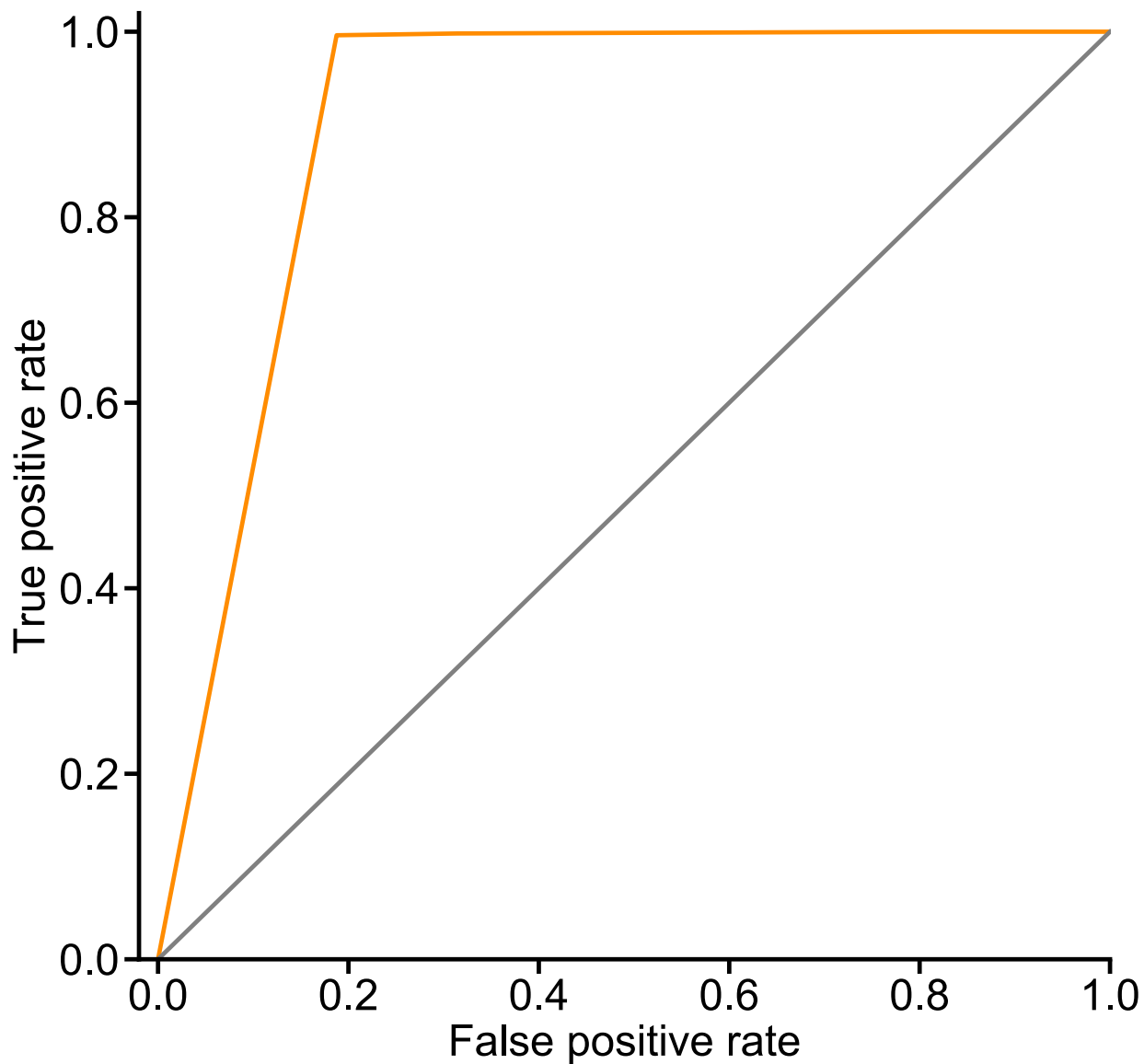


Fig. S11. AnnotSV excelled when tested on ClinVar variants that have a high degree of overlap with its cataloged variants. AnnotSV (orange) has an AUC of 0.905 when tested on 1042 ClinVar pathogenic variants and 867 ClinVar benign variants. It predicts nearly all pathogenic variants as pathogenic, along with 19% of the benign variants. The gray line represents chance.

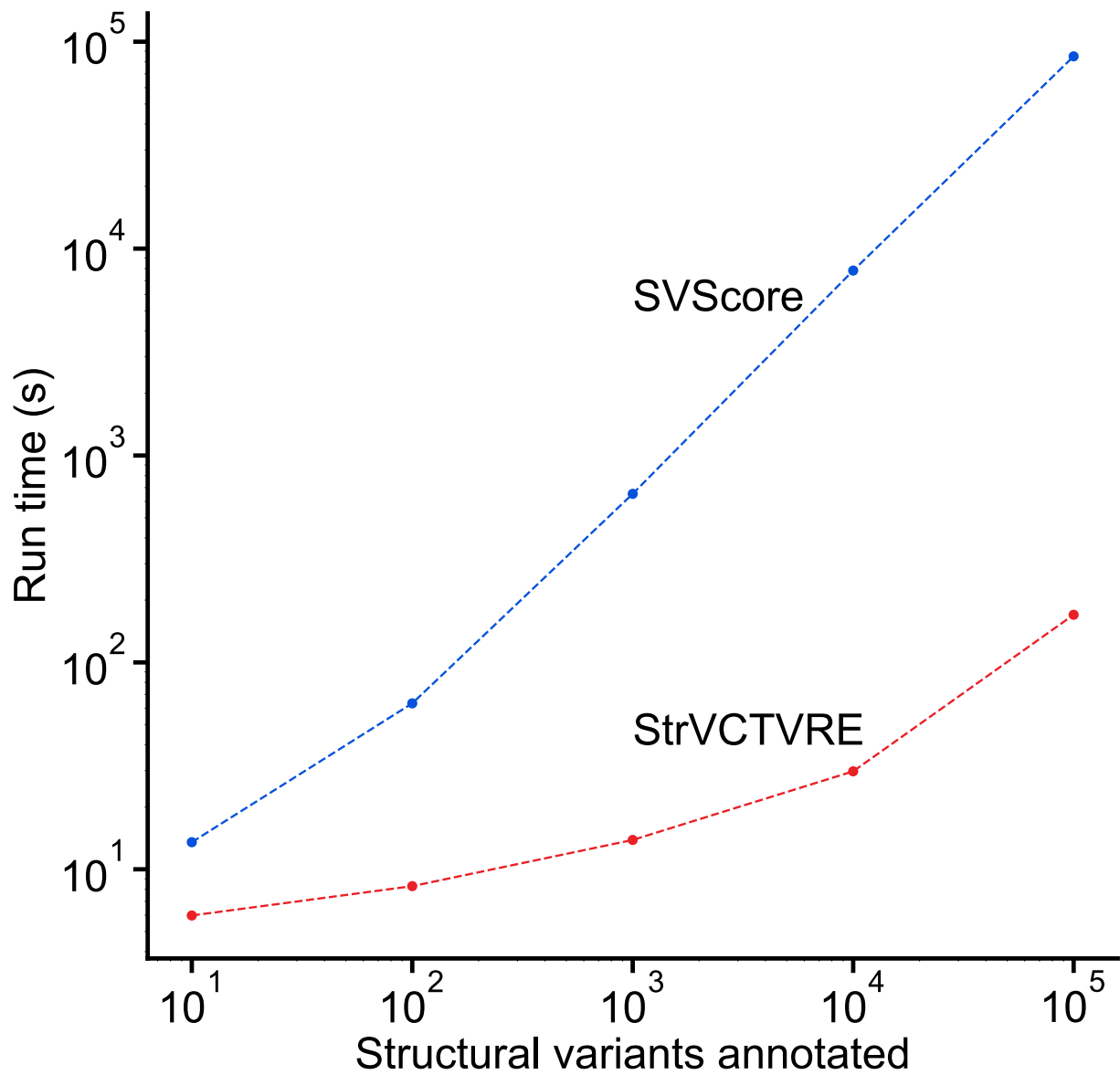


Fig. S12. Comparison of time required for StrVCTVRE and SVScore to annotate the same N structural variants. Each method was run on a 64-core Linux server with 2.6GHz Intel Xeon CPUs.

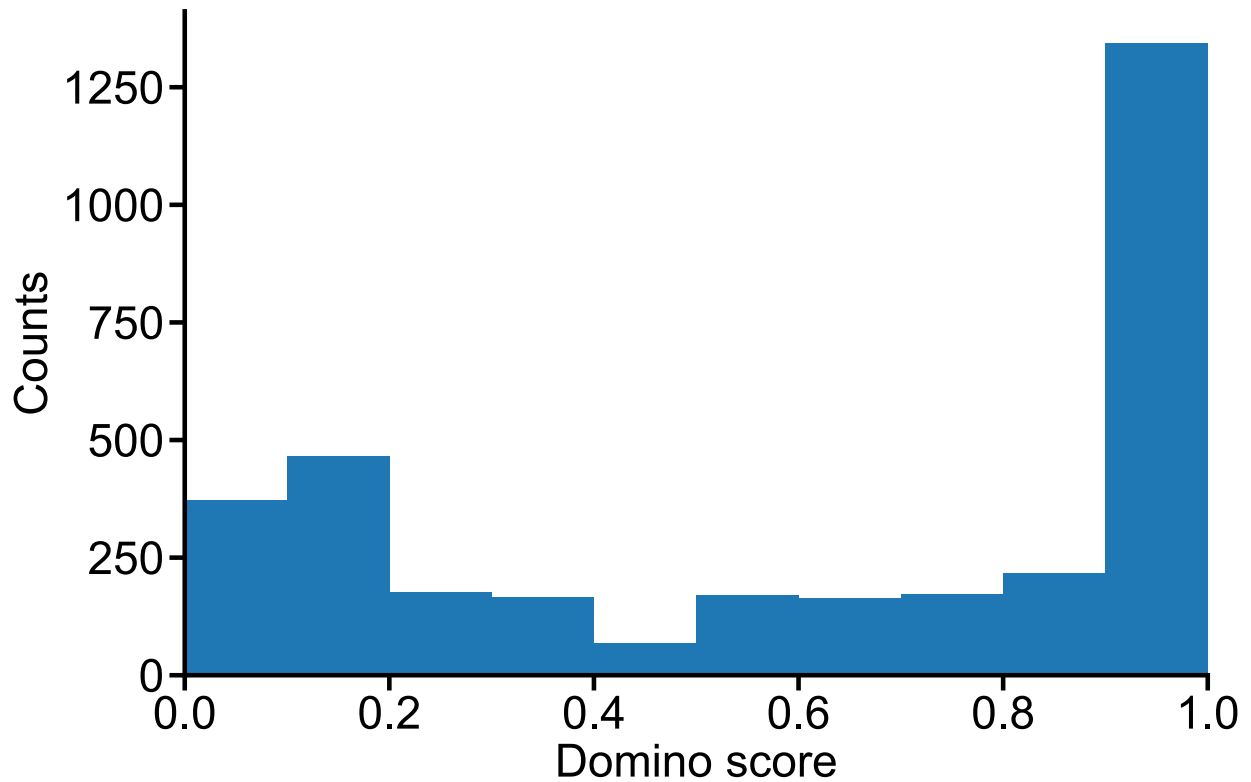


Fig. S13. Histogram of the predicted dominance of pathogenic SVs used in training. Predicted dominance was calculated for each SV by taking the maximum Domino(1) score of all genes the SV overlaps. SVs with larger Domino scores are more likely to have a dominant effect. Using a threshold of 0.5, 62% of pathogenic SVs are predicted to have a dominant effect. This analysis excludes SVs on sex chromosomes as they are not scored by Domino.

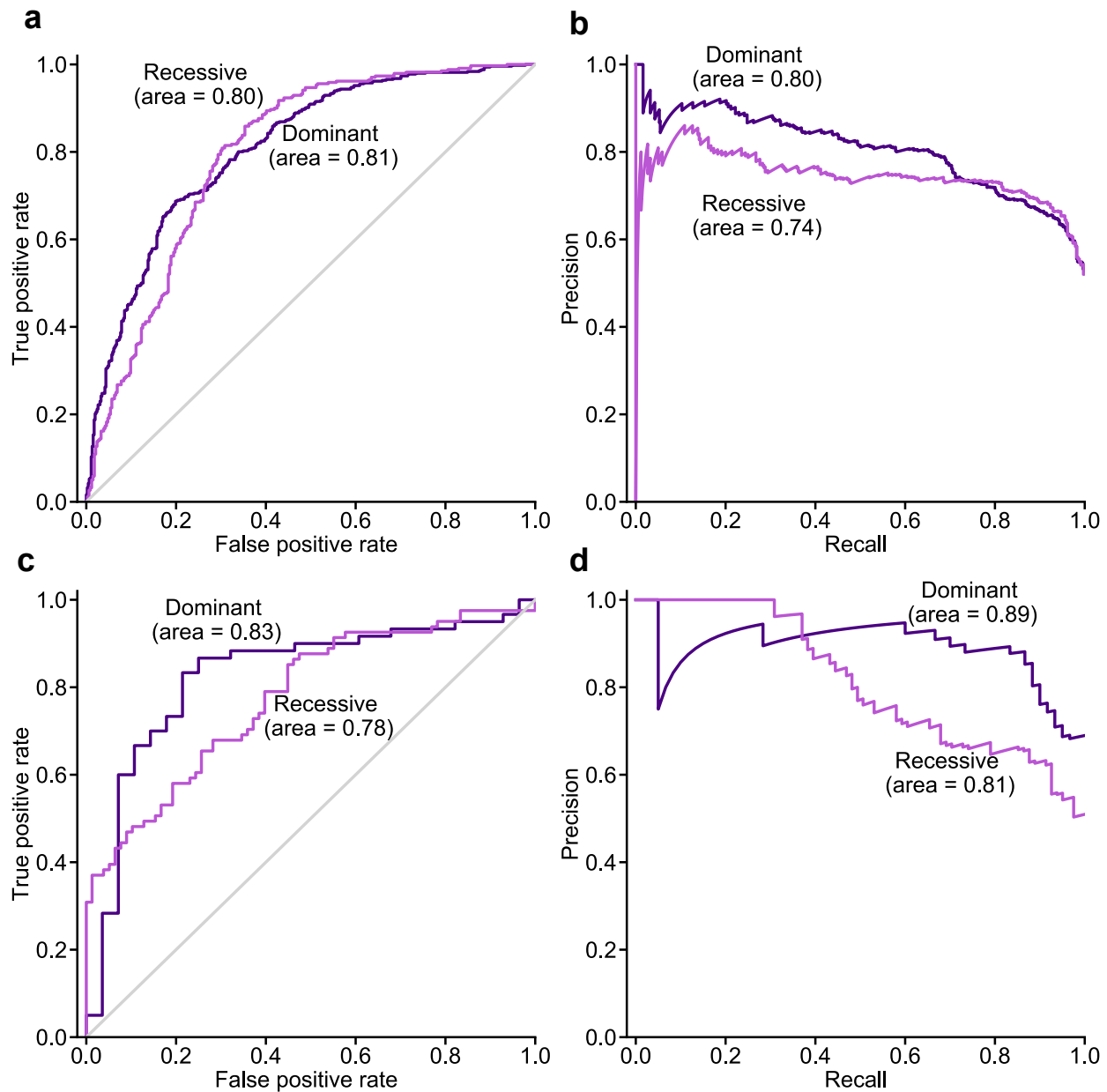


Fig. S14. ROC and Precision-Recall Curve (PRC) plots comparing StrVCTVRE's performance on recessive and dominant SVs. **a** ROC and **b** PRC plot of StrVCTVRE's performance on SVs in predicted dominant and recessive genes, using a Domino(1) threshold of 0.5 to separate dominant from recessive genes. **c** ROC and **d** PRC plot of StrVCTVRE's performance on SVs overlapping high-confidence recessive or dominant genes identified by Balick et al(2). Overall, StrVCTVRE performs similarly on SVs in both dominant and recessive genes. There may be differences in performance at low sensitivity/recall but they are not consistent across datasets.

Supplemental Tables

chrom	start	end	svtype	StrVCTVRE Score
chr1	145687715	146020436	DEL	0.768
chr1	202434559	202604719	DEL	0.673
chr2	111994154	111994818	DEL	0.569
chr2	219420423	219426832	DEL	0.511
chr6	10791578	10791945	DEL	0.276
chr6	64997485	65057825	DEL	0.634
chr6	64997485	65057825	DEL	0.634
chr10	27043329	27067384	DUP	0.699
chr13	23320540	23320858	DEL	0.3
chr15	42321688	42360212	DEL	0.502
chr15	42384386	42394439	DEL	0.619
chr16	28486250	28486550	DEL	0.571
chr17	7454200	7454618	DEL	0.567
chr17	50170063	50170449	DEL	0.471
chr19	54105500	54126715	DEL	0.678
chr19	54114345	54129468	DEL	0.753
chr19	54121739	54131817	DEL	0.847
chrX	31627576	31679684	DEL	0.835
chrX	31819878	31968612	DEL	0.847
chrX	32545062	32699391	DEL	0.823
chrX	154379176	154381621	DEL	0.601

Table S1. Subset of the 34 CMG clinical SVs used to evaluate the StrVCTVRE 90% sensitivity threshold. Due to data restrictions, the full list of SVs is not available. See 'Data and code availability' for more details.

Supplemental References

1. Quinodoz M, Royer-Bertrand B, Cisarova K, Di Gioia SA, Superti-Furga A, Rivolta C. DOMINO: using machine learning to predict genes associated with dominant disorders. *The American Journal of Human Genetics*. 2017;101(4):623-9.
2. Balick DJ, Jordan DM, Sunyaev S, Do R. Overcoming constraints on the detection of recessive selection in human genes from population frequency data. *bioRxiv*. 2021.