

SIGNIFICANCE

Structural variations (SVs), such as deletions, duplications, insertions, inversions and translocations, are among the most significant determinants of human genetic diversity to have been discovered. SVs affect far more bases than single-nucleotide polymorphisms (SNPs); thus, they can markedly affect phenotype in many ways, including modification of open reading frames, production of alternatively spliced mRNAs, alterations of transcription factor (TF) binding sites and structural gains or losses within the regulatory regions. Consortium efforts such as the 1000 Genomes Project (1000GP) estimate that a typical genome contains 2.1–2.5 thousand SVs, affecting ~20 million bases, or ~5–6 times that of SNPs. Beyond “simple” SVs, there is a growing appreciation for “complex” SVs in human genomes, which vary considerably in their architecture, ranging from small-scale insertions/deletions to complex patterns of rearrangements between distinct loci and/or even different chromosomes¹. Through the 1000GP, we found that a large fraction of SV events have much higher breakpoint complexity than previously estimated—suggesting that complex SVs, like simple SVs, are also widespread in human genomes.

Given that SVs are common, larger in size and more structurally diverse than single nucleotide variants (SNVs), they are likely to profoundly shape the regulation of many human phenotypes and disease states. Investigating SVs, and particularly complex SVs, could therefore hold the key to a deeper, more mechanistic understanding of rare and common diseases. At present, most studies do not capture the spectrum of complex SVs present in genomes, so this complexity is not adequately accounted for in disease association studies. Furthermore, the functional impact of SVs, especially in noncoding regions, has not been investigated systematically. Surmounting these issues will depend on novel computational methodologies for 1) mining whole genome sequencing datasets for SV discovery at high resolution and large scale, 2) functionally interpreting their origins and phenotypic effects, and 3) establishing associations between specific SVs and disease.

We seek support to establish The Jackson Laboratory Center for Structural Variation Analysis (JAX CSVA), to advance the overarching goals of the TOPMed program through computationally-driven discovery, functional validation and characterization of disease-associated SVs (Figure 1). We will integrate novel and powerful tools for high-resolution SV discovery and, in collaboration with the primary data-producing centers of TOPMed, use these to comprehensively profile all types of SVs, including complex SVs, from a large subset of the genomes being sequenced (Aim 1). To examine the functional impact of the identified SVs, we will integrate RNA-seq data and develop novel methodologies for functional annotation of variants and characterization of associated biological processes (Aim 2); these studies will also enable us to prioritize subsets of SVs for the association studies proposed in Aim 3. Finally, we will scale up SV detection and analysis through genotyping of all SVs detected in Aim 1 across the ~100,000 samples of TOPMed, which will provide the necessary statistical power for meaningful genotype-phenotype associations for disease-based SV association studies (Aim 3). We will be able to make inferences about human population structure and adaptation at a scale much greater than previous attempts. Our deliverables will be the largest library of validated SVs discovered in humans, together with an unprecedented platform of cloud-based pipelines for comprehensive, high-resolution and large-scale SV analysis.

INNOVATION

The originality of the JAX CSVA lies in the integration of cutting-edge computational methodologies—pioneered by the group—into a comprehensive platform for novel SV discovery, characterization and association with common human diseases. It is well known that our ability to generate large-scale genomic sequencing data is far outstripping our ability to analyze it at the scale and resolution required to make definitive functional associations. This issue is particularly relevant in the context of complex SVs, for which important details of their origin and functional effects cannot be appreciated without the proper tools for analysis at nucleotide resolution. Furthermore, the present approach combines high-resolution SV analysis balanced against the scale required for adequately powered association analyses. Our proposed detection and genotyping strategy provides higher power and resolution for investigating association between SVs spanning a large size spectrum and various phenotypes, surpassing previous standard approaches employed in current SV association studies. Briefly, the key innovations of our approach are: **1)** Development of a scalable pipeline incorporating the latest, cutting-edge SV detection and integration tools, with a focus on high-resolution classification of complex SVs. **2)** Tools for annotating SVs with functional data from coding and non-coding (nc) regions of the genome, especially through the integration of RNA-seq data. **3)** Tools for mechanistic interpretation of SVs across different classes, allowing us to make inferences about population structure and human adaptation and evolution. **4)** Association tests that integrate weighting methods for various biological

Comment [DM1]: I am not getting a full appreciation for the difference between a simple and complex SV. Can you provide once sentence about what the key difference is? Maybe move figure 4 up?

Comment [DM2]: Are you looking to create a center? I think for the small size of the grant, and the distributed nature of the groups, calling it a JAX Center may be a bit much

Comment [DM3]: Need a new figure 1

considerations, such as allele frequency and impact score, to a generalized linear model for capturing subtle association signals often missed by conventional approaches. **5) Genotyping the library of functionally and genetically relevant SVs across the entire cohort of TOPMed samples for well-powered genotype-phenotype associations in a disease context. This systematic review of complex SVs will yield the largest reference database of validated SVs to date, together with an unparalleled system for high-dimensional, high-resolution studies of SV architecture and function in health and disease.**

CENTER STRUCTURE AND ADMINISTRATION [Need a new Figure 2 for project structure]

By focusing on the discovery and analysis of SVs—a widespread form of genomic variation observed across common and Mendelian diseases—we expect to greatly enhance the ability of the TOPMed program to connect genetic variation to phenotypes of heart, lung, blood and sleep disorders. We anticipate that our analyses of TOPMed data sets, along with the analytical tools that we will make available for use in a collaborative, cloud-based environment, will be critical for maximizing the utility of the data and extracting meaningful biological insights.

Scientists participating in the JAX CSVA are leaders in SV discovery and analysis. The three PIs, Charles Lee, Ph.D., Mark Gerstein, Ph.D. and Li Ding, Ph.D., have a history of productive scientific collaboration and bring complementary experience in SV detection (Lee), functional interpretation (Gerstein) and large-scale data analysis (all), particularly association analysis (all). Each also brings significant experience in leading (1000GP SV group, Lee; modENCODE AWG, Gerstein; ENCODE networks group, Gerstein; PsychENCODE AWG, Gerstein; exRNA AWG, Gerstein) and participating in (1000GP, Lee/Gerstein/Ding; ENCODE, Gerstein; ICGC, Gerstein/Ding; KBase, Gerstein; GSP (Genome Sequencing Program), Gerstein) large-scale sequencing consortia. Under Dr. Lee's leadership, the 1000GP SV project identified SV events in ~2,500 healthy genomes and helped define the methodologies for identifying and characterizing SVs from "lower depth" (~4X) whole genome sequencing (WGS) datasets.

The proposed program will be supported by an extensive computational infrastructure at The Jackson Laboratory for Genomic Medicine, the site for the proposed JAX CSVA. Generous institutional commitments (see letter from Dr. Liu) towards development of the JAX SV Cloud will furnish the data storage and computational power needed for the formidable requirements of the project and will provide an environment for testing the cloud-readiness of the software pipelines that we will provide to the TOPMed program. The JAX CSVA will further benefit from the information technology, computational, bioinformatics, and software expertise resident at JAX.

RESEARCH STRATEGY

Specific Aim 1. Build an integrative pipeline for large-scale discovery of complex structural variation.

Rationale. Complex SV events are outside the design scope of available SV methods, yet are often of high impact and important for disease studies. To drive the discovery phase of the program, we are currently working on *fusorSV*, a framework developed by our group to discover SVs in thousands of sequenced whole genomes. *fusorSV* takes a data mining approach to SV calling by incorporating knowledge of the strengths of various existing SV callers using a truth set, and using this knowledge to perform discovery on a novel cohort of genomes. We will apply the *fusorSV* framework to a discovery cohort (10% random sampling) of individuals being sequenced by all of the Phase I/II TOPMed program projects. Using breakpoint assembly methods, we will perform *in silico* validation of the SV events and use the assembled contigs to investigate the inherent complexity prevalent at breakpoints. Ultimately, these studies will

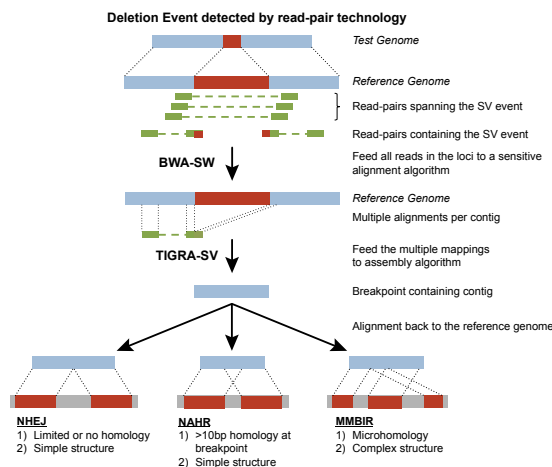


Figure 1. Breakpoint assembly for in silico validation. The top half of the figure shows a deletion SV event predicted by the readpairs spanning the event. All read pairs in the breakpoint locus are used for targeted *de novo* assembly and the resulting contig is aligned back to the genome.

Deleted: ;

Deleted: 3

deliver the largest library of validated SVs discovered in humans and allow us to make novel biological inferences at the population level and in disease-specific contexts.

Preliminary data.

A toolbox of methods for structural variation discovery. As part of the 1000GP SV project, we have provided the research community with an unprecedented set of germline SVs from more than 2,500 normal human genomes that have been sequenced at low depth and have developed a large toolbox of complementary tools and methods, including: **(i) Read depth-based tools (CNVnator).** We developed CNVnator² for copy number variant (CNV) discovery and genotyping from individual and trio-sequencing datasets. It utilizes a mean-shift approach, GC correction and bandwidth partitioning to identify a wide range of CNV events. CNVnator can detect CNVs and provide genotype information on a population level, and also detects atypical CNVs including *de novo* and multi-allelic events. **(ii) Paired end-based tools (Meerkat, Hydra-Multi, BreakDancer, Pindel).** Meerkat³, Hydra-Multi⁴ and BreakDancer⁵ cluster abnormally mapped paired-end reads to identify loci with a signature for an SV event. Meerkat remaps soft clipped and unmapped reads to generate clusters to identify breakpoints. Pindel⁶ utilizes a pattern-growth approach to detect large deletions and insertions from WGS data. These methods have individually already been successfully applied to hundreds of cancer genomes^{3,7}. **(iii) Split read alignment-based tools (SRM, SRIC and BreakSeq).** We have also developed SRM⁸ and SRIC⁹ for the high-resolution identification of SV events from WGS datasets. These tools specifically aim to provide base-pair resolution of breakpoints—an invaluable feature that enables functional interpretation of the biology of these SV events.

Breakpoint assembly tools for in silico validation. We also developed algorithms for identifying breakpoints at nucleotide resolution, thereby allowing us to validate SV breakpoints “*in silico*”. As previously described⁷, we used assembly-based methods like SGA¹⁰ or TIGRA-SV¹¹ for generating sequence contigs at breakpoints. Aligning these contigs back to the genome in the expected location and orientation validates the SV call (Figure 3). Using this method, we validated 64.8% of somatic breakpoints and 58.5% of germline control breakpoints⁷. We also developed AGE¹², which performs sequence alignment at regions flanking SVs while considering large deletion and insertion blocks, which cannot be handled by conventional sequence alignment algorithms.

Tools for complex event identification and assembly. It is now recognized that a large fraction (10–20%) of SV events are complex in nature^{7,13}. We developed PEMer¹⁴ as an initial method for identifying complex rearrangements from WGS datasets. In another study, we comprehensively characterized complex SVs from a large cohort of TCGA WGS datasets⁷ and validated them *in silico*.

Extensive complexity at structural variation breakpoints. As part of the 1000GP SV analysis team, we assessed the complexity of deletions where breakpoints had been sequenced and assembled. Consistent with the clustering analysis and the observed repeated rearrangement of duplication sites, 7.1% (1822) of these deletions intersected another deletion with different breakpoints. A larger fraction (16%) of assembled deletion sites had additional inserted sequence at deletion breakpoints. To further examine variant complexity, we grouped 1,651 deletions with at least 10 bp of additional DNA sequence between the original SV site boundaries into four broad classes (Figure 4a). The most common class, *Ins with Dup and Del*, (N=501, 30%), exhibited a recognizable duplicated sequence interval within the respective inserted sequence. Not all SVs fit neatly into the classes depicted in Figure 4a, with 214 sites forming distinct patterns exhibiting increased breakpoint complexity. Within the 1000GP sample cohort, we also found that an appreciable fraction (80%) of inversions are complex (Figure 4b), likely involving DNA replication errors^{15,16}. These results reveal the extensive complexity of SV breakpoints and highlight the importance of mining this complexity at fine resolution for interpreting the biology of SVs.

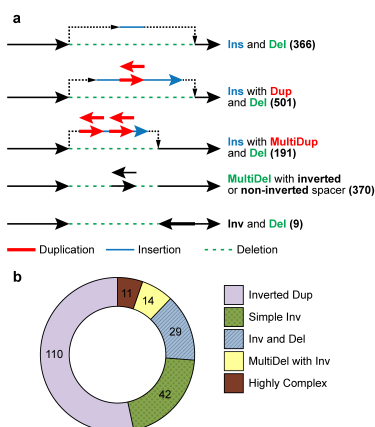
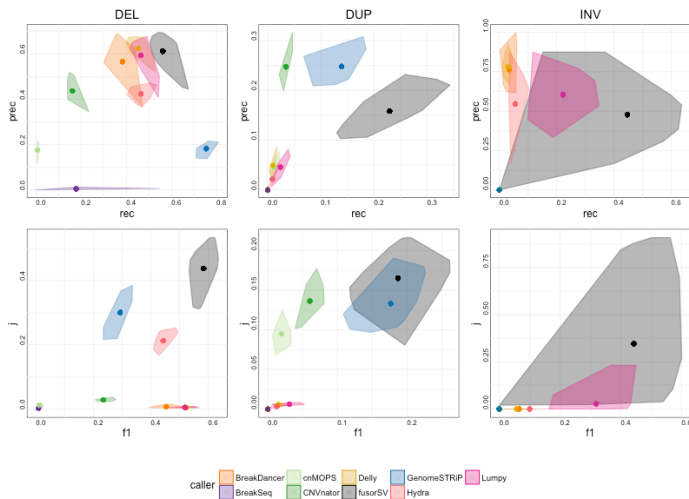


Figure 4. Structural variant complexity. a) We analyzed complexity of ~30K deletions from the 1000 GP phase 3 dataset, and characterized the events into several categories based on amount of complexity observed at the locus. b) A similar study of breakpoint complexity was performed for inversion events and revealed much higher levels of complexity than expected.



Ensemble approach to SV discovery. As noted above, *fuserSV* (manuscript in preparation) is a framework that employs a data mining approach to integrating many complementary SV callers for analyzing very large cohorts of genomes. *fuserSV* allows for germline, somatic, and *de novo* SV analysis in the cloud or on traditional high-performance compute clusters. We took deep-coverage, PCR-Free WGS data from 27 samples sequenced by the 1000Genomes Project. Using the annotated SVs from the 1000GP Phase 3, we then performed k=3 cross fold validation on this cohort, wherein we built a model using 18 samples and applied the

model to the other 9 samples for SV discovery *ab initio*. This step was repeated 1000 times with random selection for the learning samples and the test samples. Figure XXX shows the performance of *fuserSV* as compared to some of the popular SV callers that were integrated using *fuserSV*. As can be seen, *fuserSV* outperforms all the SV callers by optimizing both precision and recall on the 1000GP Phase 3 callset. Even with a harsh metric such as the Jaccard Similarity score, *fuserSV* outperforms all other SV callers for SV discovery in the test set. This Specific Aim will build on this framework and incorporate many other novel computer algorithms and improve performance.

Comment [DM4]: This figure will need a legend

Research Plan. We plan to develop new tools and extend the *fuserSV* framework to identify and classify somatic and germline SVs across WGS datasets from the various projects of the TOPMed program. The new and improved *fuserSV*, will deliver 1) integrated and comprehensive identification of a broad spectrum of SV types created by different molecular mechanisms; 2) compatibility with second- and third-generation-sequencing technologies and 3) breakpoint resolution identification based on TIGRA-SV (and other tools) and local assembly for *in silico* validation of the SV event.

Sample selection. Data storage and compute requirements preclude SV discovery on the whole TOPMed program. Based on our power calculations (Aim 3), we will select a discovery cohort of 10K individuals across all Phase I/II projects for *de novo* SV calling. This will be important to assess the applicability and efficiency of our pipeline using datasets generated from different sites. We will prioritize sample selection based on availability of orthogonal datasets (e.g., RNA-Seq, Methyl-Seq etc) and phenotypic information (e.g., blood pressure, glucose levels, BMI, etc). Clearly, having additional genomic/phenotypic data would allow us to mine better biological inferences from the SV calls.

Pipeline for population-level structural variant discovery. During phase 3 of the 1000 GP SV project, we used an ensemble of nine algorithms for SV discovery. Individual call-sets were merged into a single release through a procedure that involved re-genotyping SV genomic loci using GenomeStrip with an emphasis on genotype concordance for overlapping sites. The proposed *fuserSV* framework (Figure XXX) for SV discovery will extend this work with the following salient features: 1) MySQL database-based sample tracking of data files through the pipeline, 2) Standard steps for quality control, duplicate removal and alignment for all selected samples, 3) An ensemble of SV-calling methods including CNVnator, cnMops, BreakDancer, Pindel, Hydra-Multi, Delly, BreakSeq2, Lumpy and GenomeStrip. This ensures that a particular algorithm does not bias the discovered SV set and increases our power to detect true SV events by asking for evidence by multiple methods. The *fuserSV* framework allows us build a model using a truth set and then apply it to a novel cohort 4) Unified methods for SV genotyping and phasing using the lessons learnt from Phase 3 of the 1000GP¹⁷, 5) Validation for discovered set of SV sites using a library of known common variants and a targeted *de novo*

Comment [DM5]: Need a figure here

assembly-based approach, 6) Complex SV identification using tools for assessing breakpoints at nucleotide resolution.

The SV calling will be performed in three phases:

Phase 1—Calibration: The pipeline will use a machine-learning approach to calibrate and test the parameters of the different SV-calling methods. We will initially focus on 50 deep coverage “known truth” (KT) samples from the 1000GP SV Project¹⁷, 100 “simulated truth” (ST) samples generated using WGSim (<https://github.com/lh3/wgsim>), and 200 test cohort (TC) samples (from the TCGA consortium). These datasets all have some known true-positive SVs and will be given different weights in the eventual determination of pipeline parameters depending on the level of confidence in the associated SV set (KT>ST>TC).

Phase 2—Optimization: After calibrating our methods on the ST, KT and TC cohorts, we will expand the analysis to ~1% (~1,000) of individuals being sequenced within the TOPMed program. This cohort will be used to test for efficiency and eventual scale up in the next discovery phase. Based on the data access and compute strategies defined in TOPMed, we will explore parallelization where the tools already support this capability. The compute-intensive steps in the discovery pipeline that would be primary candidates for optimization are 1) genome alignment of raw reads, 2) clustering of aberrant reads, 3) SV validation using assembly and 4) SV integration.

Phase 3—Discovery: The optimized system will be run on 10K of the proposed 100K individuals sequenced by the various projects.

Calibration of method using known sites. Hundreds of sites across the human genome are polymorphic in a large fraction of the population^{18,19}. Phase 3 of 1000 GP SV project¹⁷ showed that a significant fraction of SVs (35%) occur at a high frequency in the population (VAF \geq 0.2%). We will create a catalog of common copy number polymorphic sites across the genome and use them as validation sites for our SV-calling methods.

Validation of SV sites using in silico assembly-based methods. We demonstrated above that SVs can be validated in silico using targeted de novo assembly-based methods (TIGRA-SV and SGA). The same methodology has been integrated into the fusorSV framework and will be used to process every discovered SV site for validation.

Complex SV identification. We will use two methods for complex SV identification. The first⁷ identifies SV clusters present in the same genomic region that have similar allele frequencies and copy number ratios. This will help select SVs that are part of the same complex SV event. The second method¹⁷ involves inspecting the mapping patterns of various parts of the assembled contig at the SV site. This would allow us to identify mislabeled SVs and SVs with more complexity than annotated by the individual SV-calling methods.

Generating a population-level reference set of SVs. We expect that several different population groups will be represented in the complete cohort of individuals being sequenced at the CCDGs and CMGs. The resulting set of validated SVs from this aim (identified from the discovery cohort) will be further stratified according to underlying population substructure. The population-specific reference set of SVs will allow us to extend the observations from 1000GP Phase 3¹⁷ and will be critical for the population and disease-level association analyses proposed in Aim 3.

Data access strategies: The JAX SV Cloud. Total storage of the discovery cohort is expected to require ~4 PB based on TCGA WGS statistics. To deal with the data footprint and computing requirements, we propose to develop the JAX SV Cloud, which will be available to all members of our teams. Our two-stage local and cloud approach is as follows:

i) *The JAX local data center.* In a traditional center, data are downloaded for analysis to local high-performance compute resources. JAX has extensive infrastructure, including an HPC cluster with 1700 cores and 1.4 PB of storage that will further expand over time (see Facilities and Resources). We can analyze the full discovery cohort by transient download and analysis of raw data with retention of only necessary results.

ii) *Cloud-based data access model.* However, it is more likely that the TOPMed program will provide access to the data using a public cloud service provider. After initial method development and analysis, we plan to disseminate methods to the broader research community using the cloud paradigm decided by the TOPMed program. JAX is currently expanding capabilities in cloud-based data analysis to address issues including access to increased compute power, co-localization of novel and reference datasets and reproducibility of analysis pipelines. JAX staff have adapted multiple pipelines for the Amazon cloud and evaluated the suitability of Amazon archival storage for genomics datasets. Dr. Ding's group has been developing GenomeVIP, a secure, HIPAA-compliant, web-driven variant discovery and annotation platform through which multiple independent analysis tools can be applied to a given dataset. As it can call upon both

Comment [DM6]: This deleted part is not pertinent to the point here, which is that there are known high-freq SVs

Comment [AM7]: Replace with a section about how 10% of individuals from each project would be selected as a discovery cohort.

Comment [AM8]: Is that what we are calling it ?

Comment [DM9]: Need to update this

local HPC and Amazon cloud resources, GenomeVIP is a tool that we may initially use to assist with variant discovery and to download results to local disks for subsequent analyses.

JAX is now evaluating commercial genomics cloud service providers (CSPs) to partner with (see letters from Seven Bridges Genomics and IBM) and is now recruiting 2 full time employees for this effort. These activities are independent of this U01 proposal. JAX will choose a platform that will allow scientists without special training to analyze their datasets through a graphical interface for both local and cloud analysis methods. This goal parallels that of the U01, namely to ensure methods developed at the data center will be stable and easily usable by the general research community.

Expected results. These studies will yield a comprehensive catalog of validated complex SVs from healthy and diseased individuals that lay the foundation for subsequent functional interpretation and association studies (Aims 2,3). They will also help answer questions about complex SV formation and population-level associations of SVs across multiple studies, thereby adding value to the TOPMed datasets. By making the fuserSV pipeline available as a community resource, we expect this work to propel future genome-level SV analyses.

Pitfalls and alternative approaches. XXX

Aim 2. Develop tools to analyze the functional impact of structural variants.

Rationale. Little is still known about the functional impact of SVs at a genome-wide level. SVs are disproportionately observed in the non-coding part of the genome; hence, comprehensive assessment of the functional impact of SVs will likely require the integration of large-scale data resources such as ENCODE, 1000GP and GTEx. We will create SV Impact (SVIM), a new analysis tool that integrates a myriad of datasets including existing annotations, allelic activity from RNA-seq and also eQTLs from RNA-seq to functionally prioritize SVs in preparation for disease association studies.

Preliminary data.

Tools for assessing functional impact of genomic variation in genes and pseudogenes. We developed Variant Annotation Tool (VAT) to annotate the impact of protein sequence mutations. VAT provides transcript-specific annotations of point mutations and, specifically, indels according to synonymous, missense, nonsense or splice-site-disrupting changes²⁵. We observed that genes tolerant of loss-of-function (LoF) mutations are under the weakest selection. In 1000GP Phase 3, we found that a typical genome contains ~150 LoF variants and discovered significant depletion of SVs (including deletions, duplications, inversions and multiallelic CNVs) in the coding sequences, untranslated regions and introns of genes compared to a random background model, implying strong purifying selection.

Tools for evaluating functional impact of variation in non-coding (nc) RNAs and regulatory regions. We developed tools to specifically analyze ncRNAs. Our incRNA pipeline combines sequence, structural and expression features to classify newly discovered transcriptionally active regions into RNA biotypes such as miRNA, snRNA, tRNA and rRNA²⁷. Our ncVar pipeline further analyzes genetic variants across biotypes and subregions of ncRNAs, e.g., showing that miRNAs with more predicted targets show higher sensitivity to mutation in the human population²⁸.

To better understand nc regulatory regions, we developed tools to analyze ChIP-Seq data to identify genomic elements and interpret their regulatory potential. PeakSeq identify regions bound by TFs and chemically modified histones^{29,30}. PeakSeq has been widely used in consortium projects such as ENCODE^{29,31}. The second generation of PeakSeq is a newly developed tool that uses multiscale decomposition to help identify enriched regions in cases where strict peaks are not apparent and robustly calls both broad and punctate peaks³⁰. Peak calls and ChIP-Seq signal data can also be used to model gene expression and annotate target genes. We have developed methods that use both supervised and unsupervised machine-learning techniques to identify these regulatory regions (such as enhancers) and predict gene expression from ChIP-Seq data³²⁻³⁵. To investigate the evolutionary importance of these regions, we have analyzed patterns of single nucleotide variation within functional nc regions, along with their coding targets^{28, 35,36}. We used metrics such as diversity and fraction of rare variants to characterize selection pressure on various classes and subclasses of functional annotations²⁸. We have also defined variants that are disruptive to a TF-binding motif in a regulatory region³¹.

Tools for helping annotate functional impact based on network. We found that functionally significant and highly conserved genes tend to be more central in various biological networks³⁷ and are positioned at the top of regulatory networks³⁶. Further studies showed relationships between selection and protein network topology (e.g., quantifying selection in hubs relative to proteins on the network periphery^{37,38}). Incorporating multiple

Comment [DM10]: Do we have these somewhere?

Comment [AM11]: Add text about how fuserSV can and should be used by the whole TOPMed program for SV discovery

Deleted: SVs account for more nucleotide variation in the human genome than SNPs and therefore are likely to be associated with many genetic diseases. However,

Deleted: their

Deleted: This proposal will catalogue the largest number of SVs so far and, more importantly,

Deleted:

Deleted: data

Deleted: annotated variants from 1,092 humans in Phase 1 of the 1000GP²⁶ and

Deleted: and cancer-causal genes under the strongest selection

Deleted: p

Deleted: s

Deleted: and allelic expression analyses

network and evolutionary properties, we developed NetSNP³⁷ to quantify the indispensability of genes. This method shows strong potential for interpreting the impact of variants involved in Mendelian diseases and in complex disorders probed by GWAS. We constructed regulatory networks for data from the ENCODE and modENCODE projects, identifying functional modules and network hierarchy³⁶. To quantify the degree of hierarchy for a given hierarchical network, we defined a metric called hierarchical score maximization (HSM³⁵).

Deleted: analyzing

FunSeq: Tools for integrated functional prioritization. We recently developed a prioritization pipeline called FunSeq^{26,40} that identifies annotations under strong selective pressure as determined using genomes from many individuals from diverse populations. FunSeq links each nc mutations to target genes and prioritizes based on scaled network connectivity. FunSeq identifies deleterious variants in many nc functional elements, including TF binding sites, enhancer elements and regions of open chromatin corresponding to DNase I hypersensitive sites, and detects their disruptiveness in TF-binding sites (both LoF and gain-of-function events).

Deleted: single-nucleotide

Mutational mechanisms of structural variants. The sequence content of SVs, especially around breakpoints, carries important information about origin and functional impact. Using datasets from 1000GP, we have studied the distinct features of SVs originating from different mechanisms^{26,42}. We performed SV mechanism annotations for the 1000GP Phase 3 deletions using BreakSeq⁴³, categorizing 29,774 deletions by their creation mechanisms. Among these, NHR proved to be the most prevalent mechanism (~73% of all categorized deletions)¹⁷. These results inform us on the molecular mechanisms underlying SV formation and also indicate differences in functional impacts of different SV types.

Deleted: We further enhanced FunSeq (FunSeq2) and identified ~100 nc candidate drivers in ~90 WGS medulloblastoma, breast and prostate cancer samples²⁶.

Deleted: For example, non-allelic homologous recombination (NAHR), is associated with active enhancers and an open chromatin environment. Our analysis also showed that micro-insertions, flanking non-homologous breakpoints, originate from late-replicating genome loci with characteristic distances from breakpoints.

Deleted: further

Tools for uniform processing of RNA-seq data. We have considerable expertise in analyzing RNA-Seq data, including experience in developing and setting up pipelines for the processing of RNA-seq data; specially for long RNA-seq data for ENCODE, long and short RNA-seq data for the PsychENCODE^[cite{26605881}] and Brainspan project as well as a custom pipeline developed for the analysis of small exRNA-seq data for the Extracellular RNA Communication Consortium (ERCC). We have already developed an efficient in-house data processing workflow for RNA-seq data that includes data organization, format conversion, and quality assessment. RSeqTools^[cite{21134889}] is a modular tool developed for the processing of RNA-seq data and generating either transcript, gene or exon level quantifications. We also developed IQSeq^[cite{22238592}] which calculates the relative and absolute abundance of contributing transcript isoforms to a gene from RNA-seq data using a fast algorithm based on the Fisher information matrix. Another tool we developed called FusionSeq^[cite{20964841}] was to detect fusion transcript in RNA-seq data, which can be important biomarker for diseases such as various types of cancer and mental diseases.

Deleted: Our lab has

Tools for allele activity and eQTL detection. We have also developed tools specifically for linking gene expression variation to genotype, including our Allele-Seq pipeline, which quantifies allele-specific gene expression by mapping reads onto a diploid personal genome built from called genetic variants, including SNPs, short indels, and structural variants^[cite{44}]. We recently applied this pipeline on a population scale to RNA-Seq data from the 1000 Genomes Project, and used this analysis to create AlleleDB, a database of genomic regions with high allelic activity^[cite{27089393}]. Our expertise in eQTLs is demonstrated in our novel study on successfully utilizing expression-variant correlations to construct predicted genotypes. These predicted genotypes were then matched with known genotypes from a given dataset in order to demonstrate how the information security of the given dataset may be compromised^[cite{26828419}].

Deleted: The software we have written for this analysis is available online (privaseq.gersteinlab.org).

Research plan. To enable identification of SVs with high functional impact, we will extend FunSeq/FunSeq2 within a new pipeline called SVIM (Structural Variation Impact)(Figure 6). We will evaluate the impact score for each SV, taking into account the functional annotation of the affected genomic region and the fraction of functional elements (i.e., genes, ncRNAs, nc regulatory elements). We will also up weight SVs based on ubiquitous activity, allelic activity and eQTLs. The impact score will also depend upon SV type (i.e., deletion, duplication, inversion or translocation).

Comment [DM12]: Aims dependency?

For a given SV belonging to a particular SV type, we will use break point resolution coordinates to estimate the fraction of bases overlapping functional elements. Based on this fraction, we will categorize SVs into three classes (touch, cut, and engulf). Each overlapping class will have a different weight ($F_{svtype, class}$). We will divide genomic elements into three categories (coding region, nc region, TF binding site) and assign relative scores to them (S_{coding} , $S_{non-coding}$, S_{TFBS}), which will vary for different SV types. Relative scores F and S will be defined for class and functional elements analogous to the FunSeq2 tool²⁶.

Deleted: five

Deleted: , allelic activity and eQTL

Deleted: , $S_{allelic}$, S_{eQTL}

$IS_{orig} = \sum_i (F_{j,k} \times S_{j,i} \times \delta_i) \times \prod_l g_l$; $IS_{norm} = \frac{IS_{orig} - IS_{random}}{\sigma_{random}}$, where i is a functional element \in {protein coding, noncoding RNA, noncoding regulatory, allelic activity, eQTL}; k is a overlapping classification \in {cut ($0.1 \leq f < 0.8$), touch ($f < 0.1$), engulf ($f \geq 0.8$)}, and f is the fraction of functional element overlapping the SV; j is the type of SV; $\delta \in \{0,1\}$; and l is a feature \in {connectivity, ubiquitous activity, allelic activity, eQTLs};

Deleted: $\prod_l g$

SVs will be assigned an impact score by taking the sum over the product between weights of overlapping classes and scores of overlapping functional elements. The score (IS_{orig}) will also be upweighted based on activity of the affected region. The upweight factor is comprised of the product of **four** factors: i.e., allelic activity, eQTLs, network connectivity and ubiquitous activity. Significance level of an Impact score (IS_{orig}) will be estimated by running 1,000 Monte Carlo simulations generated by randomly shuffling the location of SVs.

Deleted: expression

Deleted: three

Deleted: transcription

Evaluating effect of structural variants on protein-coding genes. We will further develop a protein-coding module for SVIM to substantially expand the analysis of loss of function (LoF) variants with mis-mapping, functional, evolutionary and network features. We will first identify LoFs due to whole gene deletion, as well as putative LoF-causing mutations as those that induce premature stop codons, frameshifted open reading frames, or that we predict to produce truncated proteins due to deletion of RNA splice sites or either predicted or verified changes in splicing pattern from RNA-Seq data (see above). We will quantify the confidence of these LoFs using features such as whether they are in highly duplicated regions and the number of paralogs. For functional features, we will incorporate protein structures. For evolutionary properties, we will quantify the conservation of LoF variants, as well as truncated sequences. For network features, we will quantify the distance between genes with LoF variants and known disease-causing genes.

Deleted: ui

Prioritizing non-coding transcripts from structural variant data. To prioritize the effects of SVs in ncRNAs, we will focus on overlaps with regulatory elements and other functional regions. To perform this analysis, we will define categories of RNA regions that display human population-level conservation, and combine these features to generate RNA element scores. We will mine RNA interactions between proteins (e.g., CLIP-Seq) and miRNAs (e.g., TargetScan) to create a compendium of biochemical interactions with RNA

Deleted: Finally, we will develop a machine-learning method to quantify whether LoFs will cause benign, recessive or dominant disease-causing effects. Given that most rare variants are heterozygous, developing methods to differentiate benign rare variants from disease-causing variants in terms of those that can lead to recessive or dominant disease are much needed.

⁴⁵⁻⁴⁹ We will further investigate RNA secondary structure, looking for structured regions that are highly sensitive to mutation. For these regions, we will assess deleteriousness of mutations by differences in predicted free energy or structure ensembles ⁵⁰ relative to wild type. We have found annotations of all of the above types—biochemical interactions, regulatory motifs, and structured regions—that are enriched for rare variants in the human population and will use these sensitive RNA regions to score and prioritize potential deleterious SVs in ncRNA. Large SVs will ultimately be scored based on the highest scoring subregion disrupted (or created) by the SV.

Comment [DM13]: Need to redo this reference

Comment [DM14]: Redo this reference

Prioritizing non-coding regulatory elements from structural variant data. Unlike protein-coding genes and ncRNAs, TF binding motifs are relatively small in size. Thus, we are going to analyze duplications that occur close to these motifs and analyze where these duplications lead to the breakage of existing or creation of new motifs. In the prioritization scheme, we will also penalize changes in distance between motifs and newly created motifs if they occur close to an existing TF motif. We will use TF binding nc elements by leveraging better enhancer definitions provided by the Epigenome Roadmap ⁴⁵⁻⁴⁷ and ENCODE and also include new datasets.

Deleted: leavareging

Further variant prioritization based on networks, tissue specificity, eQTLs and allelic activity. After performing annotation-based assessment of identified SVs, the following functional features will be used for prioritization.

i) Network connectivity. We will examine the network topological properties of the genomic elements affected by identified SVs. Variants disrupting regulatory elements with high connectivity—network hubs and bottlenecks—will be upweighted based on their scaled centrality scores.

ii) Ubiquitous activity. We will evaluate the impact of SVs in an epigenetic context to identify tissue-specific phenotypic effects that are strongly influenced by SVs. We will prioritize SVs impacting genes, ncRNAs, and TF binding sites active in multiple tissues.

Deleted: specificity

iii) Allelic activity. We will use our existing AlleleSeq pipeline to annotate the transcripts produced at SV regions ⁴⁴. We will use this tool to create personal diploid genomes for each TopMed individual, and then will adapt our pipeline to perform RNA-Seq quantification specifically at SV regions. We will prioritize SVs that lead to strongly allelic expression. We will also prioritize SVs that overlap our database of strongly allelic regions

throughout the genome, based on AlleleDB, our resource of such regions identified through allele-specific RNA-Seq analysis from over 300 individuals generated by the GEUVADIS consortium [cite{27089393}](#).

iv) eQTL association. We will link SVs to the genes that they affect by performing genome-wide searches for eQTLs. Relative to SNVs, large SVs may be more manageable candidates in the search for distal eQTLs. We will use a framework similar to [published earlier cite{26828419}](#) in the search for SV-induced eQTLs. SV-induced eQTLs will be identified by performing genome-wide searches for patterns in which the presence or absence of the SVs (from Aim 1) strongly correlate with the expression levels of a battery of genes throughout the genome. Specifically, we will use Matrix eQTL for eQTL identification [cite{22492648}](#). We will perform multiple testing correction and will filter the list of putative eQTLs in order to achieve a false discovery rate of less than 5%. The SV-gene expression correlations reported by Matrix eQTL will be used as the strength-of-association measures between expression levels and genotypes. Of particular interest will be those genes previously implicated in disease-associated pathways and network modules. SV-induced eQTLs with strong expression correlations that are associated with central network elements and known disease-associated genes will be upweighted.

Expected results. We expect that SVIM, a new software solution to estimate the impact scores of the SVs produced in Aim 1, will yield a prioritized set of SVs in Aim 2 that we can forward to Aim 3 (genotype and association) for further classification of their association to disease or a specific phenotype. We plan to make the prioritization results broadly available; therefore, SVIM will incorporate the impact score into a standard Variant Call Format (VCF). SVIM will be cloud-ready and will be available to the TopMed consortium through a Docker image and a Common Workflow Language (CWL) file.

Pitfalls and alternative approaches. We anticipate that the greatest pitfalls are (i) possibly an overwhelming number of SV discovered in Aim 1 and (ii) the lack of standard format and increasing number and updates of annotation datasets. In order to overcome (i), we plan to gradually process the results into specific type of SVs. SVIM will also be based on the data context to optimally prioritize from WGS datasets. The overall modularization offers a flexible framework for users to incorporate the ever-increasing amounts of genomic data to both rebuild the underlying data context and prioritize case-specific variants. In order to overcome pitfall (ii) we will make great efforts to make SVIM computationally efficient and able to support the large-scale computing proposed for this aim. To build the data context, we will [standardize](#) large-scale publicly available data resources, such as SVs from the 1000 GP⁴⁸, conservation data from Bejerano *et al.*⁴⁹ and Cooper *et al.*⁵⁰, functional genomics data from ENCODE³¹ and Roadmap Epigenomics Mapping Consortium⁵¹.

Aim 3. Scaling up to 200K samples and associating SVs with common and rare diseases.

Rationale. As many of the high-impact SVs will be relatively rare, such that conventional association tools cannot readily, robustly handle them, we will not only discover important SVs, but will develop a new association pipeline tailored for finding important SV-phenotype associations. We anticipate that building a reference database of complex structural variants in healthy individuals (Aim 1) will be essential for this goal.

Preliminary Results

- Deleted: . [1]
- Deleted: specific
- Deleted: comprehensive
- Deleted: Indels have a greater tendency to affect gene expressions relative to SNVs cite{26828419}
- Deleted: thus
- Deleted: PrivaSeq
- Deleted: the software program

Deleted: as part of a scoring scheme to be integrated into SVIM

Comment [DM15]: Can we specifically say something about making SVIM cloud-ready?

Deleted: to be

Deleted: standartize

Comment [DM16]: This creates an aim dependency, which might be a problem

Power analysis for sample selection and association. An important aspect will be selecting of a subset of the 200K samples projected to be sequenced for full SV analysis. This “discovery phase” will furnish the prototype events that will subsequently be studied in the full population by genotyping the entire sample set. Total analysis cost (e.g. downloading, storage, compute time, manual review) must be balanced against the discovery probability for events having the lowest population minor allele frequency (MAF) we wish to include. There is no general theory of discovery power currently used in SV algorithms, so we extended an existing statistical model of coverage⁵² to estimate the discovery sample size. Bernoulli probabilities for two standard SV discovery methods, split reads and discordant read pairs, can be derived using probability theory considering read length, average and variance of insert length, SV length, etc. and subsequent incorporation of a detection rule, e.g. “≥3 split or discordant reads”. Detection in each sample is binomial in the number of observations and discovery within sample set is likewise binomial in the detection and MAF probabilities.

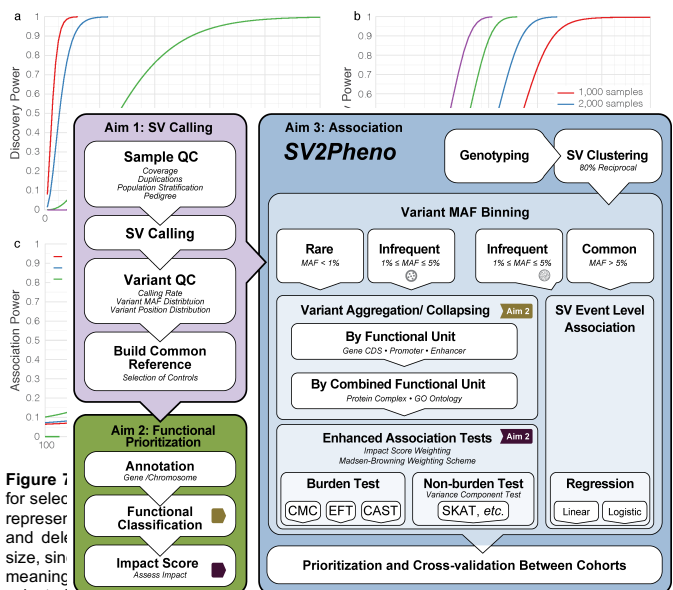
Comment [DM17]: Is this an accurate number?

Anticipated parameters for the Illumina data to be generated for this project are 30X coverage per genome, average insert size of 400bp-600bp (20% coefficient of variation), 150bp reads, event detection based on ≥3 split reads or ≥5 discordant read pairs, and observation in at least 3 samples to constitute “discovery”. The model predicts that split-read detection will predominate for simple SVs, as well as for complex events in which one sequence is replaced by another. Because split-reads depend only upon local alignment, power is essentially independent of the size of events (unlike for discordant read pairs), meaning it is primarily a function of sample size and MAF. **Figure 7a** shows power at MAF ≥ 0.1% is essentially 100% for 10K samples. It drops rapidly for lower MAFs, whose events are unlikely to be discovered in this study. Mosaicism is a potentially confounding factor, for example in blood samples where an event is not present in all cells. **Figure 7b** shows that power is not significantly impacted for 10K samples until mosaicism is quite significant.

Comment [DM18]: Double check that data is all Illumina

The second aspect of “power” is variant-disease association. The issues are well-known⁵³, enabling the following “baseline” estimates of association power. General consensus⁵³ recommends “collapsing” variants for low MAF in order to aggregate effects for increasing power. Analysis of the widely-used Li & Leal method for 10 collapsed variants at 4:1 risk ratio (**Figure 7c**) shows that groupings of 1% MAF variants having high (~50%) penetrance will require

20K-30K samples for 80% power when Bonferroni-corrected. Power drops rapidly for lower MAF, penetrance, risk ratio, and sample size. Although it is not yet known how the 200K samples will be divided over various studies, it is instructive to examine the scenario of 10K cases/10K controls (**Figure 7d**). Variants around 2% MAF should have ≥90% association power for penetrances ≥50%, while variants regardless of MAF having penetrances <25% will likely remain ambiguous, as will variants from phenotypes having substantially smaller sample allotments. It is likely we will discover more variants than those for which solid associations can be established.



Association pipeline implementation and experience in discovering significant associations. We have developed a prototype pipeline incorporating extensive sample and variant level quality control (e.g. coverage, variant frequency and distribution), population stratification, pedigree segregation

Figure 8. SV2Pheno Association Analysis Pipeline. The overall work flow (even n) includes QC, population stratification from Aim 1, functional classification and 50%, at impact score generation from Aim 2 and single event test and burden analysis 4:1 risk ratio from Aim 3. n and 200,000, collapsing strategy, a) and b) are similar to c) for 10K cases/10K controls, with other parameters the same as in c).

etc. for population/family-based association analysis. It sports popular aggregation tests, including burden tests such as the Combined Multivariate Collapsing (CMC)⁵³, Exclusive Frequency Test (EFT)⁵⁴, Total Frequency Test (TFT)⁵⁴, and Cohort Allele Sum Test (CAST)⁵⁵, and variant component tests such as the Sequence Kernel Association Test (SKAT)⁵⁶. We have already used it to discover associations by tailoring it to hypothesized genetic architectures of individual diseases. For example, assuming tumor suppressors are enriched for rare deleterious truncations, we grouped events by gene and used TFT to associate 13 genes with germline susceptibility in a >4,000 case cancer cohort.

Comment [DM19]: Reference?

Research Plan. SVs are characterized by size, type, penetrance, and multiple alleles. We plan to genotype all SVs detected in 10K discovery samples (**Aim 1**) across all ~200K samples to be sequenced by TOPMed centers to obtain sufficient statistical power for genotype-phenotype association. A critical step for association analysis of SVs is meaningful classification/annotation. By building on infrastructure and tools mentioned above, we will develop a new pipeline called “SV2Pheno” to infer SV-phenotype associations (**Fig. 8**). It will use the impact scores for each SV (**Aim 2**) for integrated analysis of SNVs, indels, and SV.

Genotyping of SVs detected in the discovery set across the entire sample set. Genotyping and annotation of discovered SVs in the whole population will allow accurate determination of prevalence and allele frequencies and, importantly, increase association analysis power. This process will use BreakSeq⁴³ to build a library of validated and assembled SV breakpoints for genotyping individual genomes. For imprecise SVs, a combined read-pair/read-depth approach using GenomeStrip⁵⁷ will do population level genotyping. Conventional genotyping involves assembly of both reference and alternate sequence contigs, which are used as targets for mapping all reads present in the sample. However, given an expected data footprint of ~50PB for the full sample set, the traditional “bring data to the computing tools” approach will be upended to “bring compute tools to the data”. We shall build on tools such as Sambamba (bam slicer function)⁵⁸, TIGRA-SV assembler and Pindel. This will reduce the footprint to a fraction of the original and enable the methods to work in the cloud and access data over a secure network.

Develop SV2Pheno pipeline including improved burden tests considering impact score and annotation classification of various complex structure variants. We envision substantial extension of this pipeline in two major ways to address the ambitious goals of this proposal: 1) We plan to hybridize the pipeline with more recent methods that better account for non-contributing variants⁵⁹. Likewise, annotation and functional prediction can help identify irrelevant variants, which can subsequently be removed from analysis. The pipeline will also process the information from the ENCODE & Epigenetics Roadmap analysis. 2) Variants are known to be associated with various diseases⁶⁰⁻⁶², but almost certainly contribute non-uniformly; assigning appropriate weights will be necessary to wring-out maximum power. Aggregation tests can be expressed in general by the linear regression equation $Y = \alpha + \beta \cdot \sum w_i g_i + \epsilon$, where (left-to-right) is observed trait, intercept, collective effect coefficient, weight of variant i , tally of variant i (0, 1, or 2), and normally-distributed error residual. Assignment of weights will be based on a novel combination of 4 considerations: the Madsen-Browning equation⁶³ to account for allele frequency, consideration of “direction” (negative association) using e.g. aspects of the Pan-Shen approach⁶⁴, incorporation of our impact score (**Aim 2**) to account for biological strength, and RNA-seq data. The last aspect will weight expression impact, but must be implemented carefully because of variations in sample quality. Here, we will apply the method of Liu et al. which essentially adds an extra calculational adjustment to modulate contribution of higher-variability samples. In principle, this more sophisticated approach should capture signals that have been too subtle for earlier-generation tests⁶⁵.

Comment [MW20]: R Liu et al (2015) Why wichtig? Modeling sample and observational level variability improves power in RNA-seq analysis. NAR. Doi 10.1093/nar/gkv412

We are mindful that controls for each association analysis should be carefully matched with cases; paying close attention to population structure, sample coverage, etc. When sample size is fixed, an even case-control split offers maximal power. However, it is likely that the TOPMed program will furnish potentially many more controls and this increases power. For such diseases, we will check the available literature for any known underlying genetic commonalities and choose extra controls in light of relevant covariates (e.g. age or smoking status). Since we anticipate that a high fraction of SVs will reside in non-coding regions, we will aggregate variants using a hierarchical approach based on three levels:

Level 1. Prototypical Event level association analysis. As the precise genomic region for a given SV may vary across samples, we will represent each set of similar SV events as a single prototypical SV event. The criterion constituting such events is given by the “80% reciprocal overlap” rule¹¹. For large insertions and inter-chromosomal translations, we will require the breakpoints to be within 1kb of one another. We will then assess the significance of the associations using impact scores generated in Aim 2.

Level 2. Functional Unit (Gene CDS/promoter/enhancer) level association analysis. We annotate the prototypical SV events from Level 1 to identify any specific transcriptional regions (e.g., exons/CDS and cis-

regulatory elements such as insulators, enhancers, and promoters) and gene(s). SVs in a given gene will be grouped as a single, effective functional unit (**Figure 8**). We will then perform an association analysis on these functional units. In cases where multiple SV events may be affiliated with a given functional unit, we will develop a weighting scheme to combine the impact scores of the contributing SVs. This approach may potentially reveal novel connections between non-coding functional regions and phenotypes.

Level 3. Combined Functional Unit level analysis. We will annotate the functional units in the previous step to identify any known affiliated higher-order units (e.g., protein complexes and gene pathways) by recruiting various resources, including databases relating to gene-phenotype relationships (e.g., OMIM), gene pathways (e.g., KEGG, Reactome), gene ontology (e.g., GO database). The SVs affecting a given higher-order unit will be grouped as a single super-unit. We will again perform association analysis, considering the SV impact scores (**Aim 2**). This approach has the potential to discover novel combinations of SV-containing functional units.

We will apply this tiered approach and association analysis (**Figure 8**) to analyze all genotyped samples passing our extensive coverage and variant calling QC from various cohorts to identify promising candidate SVs associated with specific phenotype.

Integrate various types of variants for association analysis. The most powerful analysis will come by combining information from SNVs, indels, and SVs for association analysis. Traditionally, weights in burden tests account for variants with different MAFs, but favoring those having lower MAFs^{56,63}. Bioinformatic information, such as PolyPhen scores for SNVs, and SV impact scores from Aim 2 will inform these weights. To the best of our knowledge, no previous approaches have aggregated variants of different types. Here, we propose two methods for such integration: 1) We hypothesize that SVs would have stronger functional impacts than missense SNVs, on average, and we will develop a weighing scheme based on the size and genetic architecture of various variant types using the framework of previous weighting schemes. SNV/indel/SV will be jointly calculated in a single burden analysis; 2) We hypothesize that alterations from functional regions, regardless of size, contribute to phenotype. Therefore, alternatively, we plan use SNV/indel and SV for independent burden analyses and combine the P-values from these independent tests.

Association between SNVs/indels and # of SVs. Under the null hypothesis that variation occurs randomly, it should be possible to correlate the numbers of SNVs/indels versus number of SVs, the slope being indicative of differences in rates of occurrence, and also to check such correlation against established rates. We will perform association analysis for individual outlier cases in which SV census is significantly lower or higher than expected. It is possible that such outliers might harbor common germline alterations leading to genomic instability by affecting DNA repair pathways.

Expected results. This aim will culminate in the JAX CSVA **SV2Pheno** association pipeline and its associated/support tools for systematically discovering SVs associated with specific phenotype/disease. We expect to have increased statistical power to discover rare, novel SVs associated with phenotypes previously missed due to smaller sample size. We further anticipate revealing genetic changes associated with increased frequency of SVs genome-wide. The initial version of **SV2Pheno** will be distributed for broader community use and cloud distribution.

Pitfalls and alternative approaches. Our preliminary analysis indicates that we are well powered to detect SVs with MAFs around 0.5% to 1% using 10,000 cases. Although it is very likely that we will discover more SVs than we can establish associations for (discussed above), there are still some issues of selection. There are several strategies for selecting datasets for initial discovery: 1) from one homogenous cohort; 2) from one CCDG center across multiple cohorts; 3) from multiple cohorts generated by multiple TOPMed centers. Regardless of choice, we will maintain high standards regarding coverage, read length, insert size, mapping rate, % mismatch etc. to ensure accurate, representative detection of SVs across populations. To reduce the number of hypotheses to be tested, we can alternatively focus on SVs from regions indicated to have association with phenotype from the study of SNV/indel. The weighting methods discussed above for may require tuning and we will use known disease associated SVs as positive controls for the calibration.

REFERENCES.

- 1 Quinlan, A. R. & Hall, I. M. Characterizing complex structural variation in germline and somatic genomes. *Trends Genet* **28**, 43-53, doi:10.1016/j.tig.2011.10.002 (2012).
- 2 Abyzov, A., Urban, A. E., Snyder, M. & Gerstein, M. CNVnator: an approach to discover, genotype, and characterize typical and atypical CNVs from family and population genome sequencing. *Genome research* **21**, 974-984, doi:10.1101/gr.114876.110 (2011).
- 3 Yang, L. *et al.* Diverse mechanisms of somatic structural variations in human cancer genomes. *Cell* **153**, 919-929, doi:10.1016/j.cell.2013.04.010 (2013).
- 4 Lindberg, M. R., Hall, I. M. & Quinlan, A. R. Population-based structural variation discovery with Hydra-Multi. *Bioinformatics* **31**, 1286-1289, doi:10.1093/bioinformatics/btu771 (2015).
- 5 Fan, X., Abbott, T. E., Larson, D. & Chen, K. BreakDancer - Identification of Genomic Structural Variation from Paired-End Read Mapping. *Curr Protoc Bioinformatics* **2014**, doi:10.1002/0471250953.bi1506s45 (2014).
- 6 Ye, K., Schulz, M. H., Long, Q., Apweiler, R. & Ning, Z. Pindel: a pattern growth approach to detect break points of large deletions and medium sized insertions from paired-end short reads. *Bioinformatics* **25**, 2865-2871, doi:10.1093/bioinformatics/btp394 (2009).
- 7 Malhotra, A. *et al.* Breakpoint profiling of 64 cancer genomes reveals numerous complex rearrangements spawned by homology-independent mechanisms. *Genome research* **23**, 762-776, doi:10.1101/gr.143677.112 (2013).
- 8 Malhotra, A. *et al.* Ploidy-Seq: inferring mutational chronology by sequencing polyploid tumor subpopulations. *Genome Med* **7**, 6, doi:10.1186/s13073-015-0127-5 (2015).
- 9 Zhang, Z. D. *et al.* Identification of genomic indels and structural variations using split reads. *BMC genomics* **12**, 375, doi:10.1186/1471-2164-12-375 (2011).
- 10 Simpson, J. T. & Durbin, R. Efficient de novo assembly of large genomes using compressed data structures. *Genome research* **22**, 549-556, doi:10.1101/gr.126953.111 (2012).
- 11 Chen, K. *et al.* TIGRA: a targeted iterative graph routing assembler for breakpoint assembly. *Genome research* **24**, 310-317, doi:10.1101/gr.162883.113 (2014).
- 12 Abyzov, A. & Gerstein, M. AGE: defining breakpoints of genomic structural variants at single-nucleotide resolution, through optimal alignments with gap excision. *Bioinformatics* **27**, 595-603, doi:10.1093/bioinformatics/btq713 (2011).
- 13 Weckselblatt, B. & Rudd, M. K. Human structural variation: mechanisms of chromosome rearrangements. *Trends Genet*, doi:10.1016/j.tig.2015.05.010 (2015).
- 14 Korb, J. O. *et al.* PEMer: a computational framework with simulation-based error models for inferring genomic structural variants from massive paired-end sequencing data. *Genome Biol* **10**, R23, doi:10.1186/gb-2009-10-2-r23 (2009).
- 15 Hastings, P. J., Lupski, J. R., Rosenberg, S. M. & Ira, G. Mechanisms of change in gene copy number. *Nat Rev Genet* **10**, 551-564, doi:10.1038/nrg2593 (2009).
- 16 Hastings, P. J., Ira, G. & Lupski, J. R. A microhomology-mediated break-induced replication model for the origin of human copy number variation. *PLoS genetics* **5**, e1000327, doi:10.1371/journal.pgen.1000327 (2009).
- 17 Sudmant, P. H. An integrated map of structural variation in 2,504 human genomes. *Nature Accepted, in print* (2015).
- 18 Wong, K. K. *et al.* A comprehensive analysis of common copy-number variations in the human genome. *American journal of human genetics* **80**, 91-104, doi:10.1086/510560 (2007).
- 19 Bailey, J. A., Kidd, J. M. & Eichler, E. E. Human copy number polymorphic genes. *Cytogenet Genome Res* **123**, 234-243, doi:10.1159/000184713 (2008).
- 20 Iskow, R. C., Gokcumen, O. & Lee, C. Exploring the role of copy number variants in human adaptation. *Trends Genet* **28**, 245-257, doi:10.1016/j.tig.2012.03.002 (2012).

- 21 Hehir-Kwa, J. Y., Pfundt, R., Veltman, J. A. & de Leeuw, N. Pathogenic or not? Assessing the clinical relevance of copy number variants. *Clin Genet* **84**, 415-421, doi:10.1111/cge.12242 (2013).
- 22 Almal, S. H. & Padh, H. Implications of gene copy-number variation in health and diseases. *J Hum Genet* **57**, 6-13, doi:10.1038/jhg.2011.108 (2012).
- 23 Drummond-Borg, M., Deeb, S. S. & Motulsky, A. G. Molecular patterns of X chromosome-linked color vision genes among 134 men of European ancestry. *Proc Natl Acad Sci U S A* **86**, 983-987 (1989).
- 24 Vollrath, D., Nathans, J. & Davis, R. W. Tandem array of human visual pigment genes at Xq28. *Science* **240**, 1669-1672 (1988).
- 25 Habegger, L. *et al.* VAT: a computational framework to functionally annotate variants in personal genomes within a cloud-computing environment. *Bioinformatics* **28**, 2267-2269, doi:10.1093/bioinformatics/bts368 (2012).
- 26 Khurana, E. *et al.* Integrative annotation of variants from 1092 humans: application to cancer genomics. *Science* **342**, 1235587, doi:10.1126/science.1235587 (2013).
- 27 Lu, Z. J. *et al.* Prediction and characterization of noncoding RNAs in *C. elegans* by integrating conservation, secondary structure, and high-throughput sequencing and array data. *Genome research* **21**, 276-285, doi:10.1101/gr.110189.110 (2011).
- 28 Mu, X. J., Lu, Z. J., Kong, Y., Lam, H. Y. & Gerstein, M. B. Analysis of genomic variation in non-coding elements using population-scale sequencing data from the 1000 Genomes Project. *Nucleic Acids Res* **39**, 7058-7076, doi:10.1093/nar/gkr342 (2011).
- 29 Rozowsky, J. *et al.* PeakSeq enables systematic scoring of ChIP-seq experiments relative to controls. *Nature biotechnology* **27**, 66-75, doi:10.1038/nbt.1518 (2009).
- 30 Harmanci, A., Rozowsky, J. & Gerstein, M. MUSIC: identification of enriched regions in ChIP-Seq experiments using a mappability-corrected multiscale signal processing framework. *Genome Biol* **15**, 474, doi:10.1186/s13059-014-0474-3 (2014).
- 31 Consortium, E. P. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57-74, doi:10.1038/nature11247 (2012).
- 32 Cheng, C. *et al.* A statistical framework for modeling gene expression using chromatin features and application to modENCODE datasets. *Genome Biol* **12**, R15, doi:10.1186/gb-2011-12-2-r15 (2011).
- 33 Cheng, C. *et al.* Understanding transcriptional regulation by integrative analysis of transcription factor binding data. *Genome research* **22**, 1658-1667, doi:10.1101/gr.136838.111 (2012).
- 34 Gerstein, M. B. *et al.* Comparative analysis of the transcriptome across distant species. *Nature* **512**, 445-448, doi:10.1038/nature13424 (2014).
- 35 Yip, K. Y. *et al.* Classification of human genomic regions based on experimentally determined binding sites of more than 100 transcription-related factors. *Genome Biol* **13**, R48, doi:10.1186/gb-2012-13-9-r48 (2012).
- 36 Gerstein, M. B. *et al.* Architecture of the human regulatory network derived from ENCODE data. *Nature* **489**, 91-100, doi:10.1038/nature11245 (2012).
- 37 Khurana, E., Fu, Y., Chen, J. & Gerstein, M. Interpretation of genomic variants using a unified biological network approach. *PLoS Comput Biol* **9**, e1002886, doi:10.1371/journal.pcbi.1002886 (2013).
- 38 Kim, P. M., Korbelt, J. O. & Gerstein, M. B. Positive selection at the protein network periphery: evaluation in terms of structural constraints and cellular context. *Proc Natl Acad Sci U S A* **104**, 20274-20279, doi:10.1073/pnas.0710183104 (2007).
- 39 Cheng, C. *et al.* An approach for determining and measuring network hierarchy applied to comparing the phosphorylome and the regulome. *Genome Biol* **16**, 63, doi:10.1186/s13059-015-0624-2 (2015).
- 40 Fu, Y. *et al.* FunSeq2: a framework for prioritizing noncoding regulatory variants in cancer. *Genome Biol* **15**, 480, doi:10.1186/s13059-014-0480-5 (2014).
- 41 Dees, N. D. *et al.* MuSiC: identifying mutational significance in cancer genomes. *Genome research* **22**, 1589-1598, doi:10.1101/gr.134635.111 (2012).

- 42 Abyzov, A. *et al.* Analysis of deletion breakpoints from 1,092 humans reveals details of mutation mechanisms. *Nat Commun* **6**, 7256, doi:10.1038/ncomms8256 (2015).
- 43 Lam, H. Y. *et al.* Nucleotide-resolution analysis of structural variants using BreakSeq and a breakpoint library. *Nature biotechnology* **28**, 47-55, doi:10.1038/nbt.1600 (2010).
- 44 Rozowsky, J. *et al.* AlleleSeq: analysis of allele-specific expression and binding in a network framework. *Mol Syst Biol* **7**, 522, doi:10.1038/msb.2011.54 (2011).
- 45 Roadmap Epigenomics, C. *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317-330, doi:10.1038/nature14248 (2015).
- 46 Ziller, M. J. *et al.* Dissecting neural differentiation regulatory networks through epigenetic footprinting. *Nature* **518**, 355-359, doi:10.1038/nature13990 (2015).
- 47 Leung, D. *et al.* Integrative analysis of haplotype-resolved epigenomes across human tissues. *Nature* **518**, 350-354, doi:10.1038/nature14217 (2015).
- 48 Genomes Project, C. *et al.* An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**, 56-65, doi:10.1038/nature11632 (2012).
- 49 Bejerano, G. *et al.* Ultraconserved elements in the human genome. *Science* **304**, 1321-1325, doi:10.1126/science.1098119 (2004).
- 50 Cooper, G. M. *et al.* Distribution and intensity of constraint in mammalian genomic sequence. *Genome research* **15**, 901-913, doi:10.1101/gr.3577405 (2005).
- 51 Bernstein, B. E. *et al.* The NIH Roadmap Epigenomics Mapping Consortium. *Nature biotechnology* **28**, 1045-1048, doi:10.1038/nbt1010-1045 (2010).
- 52 Wendl, M. C. & Wilson, R. K. Statistical aspects of discerning indel-type structural variation via DNA sequence alignment. *BMC genomics* **10**, 359, doi:10.1186/1471-2164-10-359 (2009).
- 53 Li, B. & Leal, S. M. Methods for detecting associations with rare variants for common diseases: application to analysis of sequence data. *American journal of human genetics* **83**, 311-321, doi:10.1016/j.ajhg.2008.06.024 (2008).
- 54 Cohen, J. C. *et al.* Multiple rare alleles contribute to low plasma levels of HDL cholesterol. *Science* **305**, 869-872, doi:10.1126/science.1099870 (2004).
- 55 Morgenthaler, S. & Thilly, W. G. A strategy to discover genes that carry multi-allelic or mono-allelic risk for common diseases: a cohort allelic sums test (CAST). *Mutation research* **615**, 28-56, doi:10.1016/j.mrfmmm.2006.09.003 (2007).
- 56 Wu, M. C. *et al.* Rare-variant association testing for sequencing data with the sequence kernel association test. *American journal of human genetics* **89**, 82-93, doi:10.1016/j.ajhg.2011.05.029 (2011).
- 57 Handsaker, R. E., Korn, J. M., Nemes, J. & McCarroll, S. A. Discovery and genotyping of genome structural polymorphism by sequencing on a population scale. *Nature genetics* **43**, 269-276, doi:10.1038/ng.768 (2011).
- 58 Tarasov, A., Vilella, A. J., Cuppen, E., Nijman, I. J. & Prins, P. Sambamba: fast processing of NGS alignment formats. *Bioinformatics* **31**, 2032-2034, doi:10.1093/bioinformatics/btv098 (2015).
- 59 Lee, S. *et al.* Optimal unified approach for rare-variant association testing with application to small-sample case-control whole-exome sequencing studies. *American journal of human genetics* **91**, 224-237, doi:10.1016/j.ajhg.2012.06.007 (2012).
- 60 Sebat, J. *et al.* Strong association of de novo copy number mutations with autism. *Science* **316**, 445-449, doi:10.1126/science.1138659 (2007).
- 61 Walsh, T. *et al.* Rare structural variants disrupt multiple genes in neurodevelopmental pathways in schizophrenia. *Science* **320**, 539-543, doi:10.1126/science.1155174 (2008).
- 62 Cooper, G. M. *et al.* A copy number variation morbidity map of developmental delay. *Nature genetics* **43**, 838-846, doi:10.1038/ng.909 (2011).
- 63 Madsen, B. E. & Browning, S. R. A groupwise association test for rare mutations using a weighted sum statistic. *PLoS genetics* **5**, e1000384, doi:10.1371/journal.pgen.1000384 (2009).
- 64 Pan, W. & Shen, X. Adaptive tests for association analysis of rare variants. *Genetic epidemiology* **35**, 381-388, doi:10.1002/gepi.20586 (2011).

- 65 Lee, S., Abecasis, G. R., Boehnke, M. & Lin, X. Rare-variant association analysis: study designs and statistical tests. *American journal of human genetics* **95**, 5-23, doi:10.1016/j.ajhg.2014.06.009 (2014).
- 66 Schnabel, Z. E. The Estimation of Total Fish Population of a Lake. *American Mathematical Monthly* **45**, 348-352 (1938).

Integration of SV annotation and RNA-seq data