# A flexible, high-throughput software pipeline for identifying candidate variants and their effects on protein structure, starting from NGS data or VCF files

Katie Maxfield, Matthew Keyser, Steven J. Darnell, Matt Larson, Erik Edlund, Richard Nelson, Schuyler Baldwin, Dan Nash and Tim Durfee
DNASTAR, Inc., Madison, Wisconsin, USA

## Abstract

Advances in DNA sequencing have made identifying genetic variation from any individual routine, allowing researchers to analyze variation across populations or cohorts for candidate mutations causing diseases or traits of interest. To address the significant bioinformatic, data management and computing resource challenges inherent to this analysis, we have developed an integrated software solution for each step in the pipeline.

First, variant call profiles for each sample are either determined directly from raw NGS sequencing data or extracted from existing VCF files. For NGS data processing, reads from each sample are aligned to a reference genome using a fast, non-memory bound assembler. Gapped alignments are analyzed in-stream using a Bayesian modeled probabilistic variant caller to produce annotated single nucleotide variant (SNV) and small indel calls. To leverage existing variant information in VCF format, we have also developed a VCF annotation tool which maps the variants onto a corresponding annotated reference sequence and "decorates" each variant with information such as the affected gene(s) and impact on protein encoding regions and/or splice sites.

Next, annotated variant profiles from each sample, supplemented with their allele frequencies as well as predictions about functional impact and pathogenicity, are automatically combined into a single project for analysis. Various filtering and statistical methods then enable candidate genes and/or variants of interest to be readily identified. Further, for genes with known 3D protein structures, the effect of candidate missense mutations on the structure can be automatically predicted and analyzed using our molecular structure visualization module.

As a demonstration of the pipeline, we present results from our reanalysis of 96 targeted resequencing samples from a Chinese cohort with lung squamous cell carcinomas (LSCC, Li et al., Sci Rep. 2015. 5:14237). We will show how the software can be used to easily identify unique mutations in numerous samples across the cohort which all lead to nonsense and frameshift mutations in the TP53 tumor suppressor gene. We will also show how filtering on the functional impact and pathogenicity predictions can be used to identify likely deleterious TP53 missense mutations in other members of the cohort and how in one of those cases, the predicted protein structure change indicates that the DNA binding activity of the protein is likely directly affected.
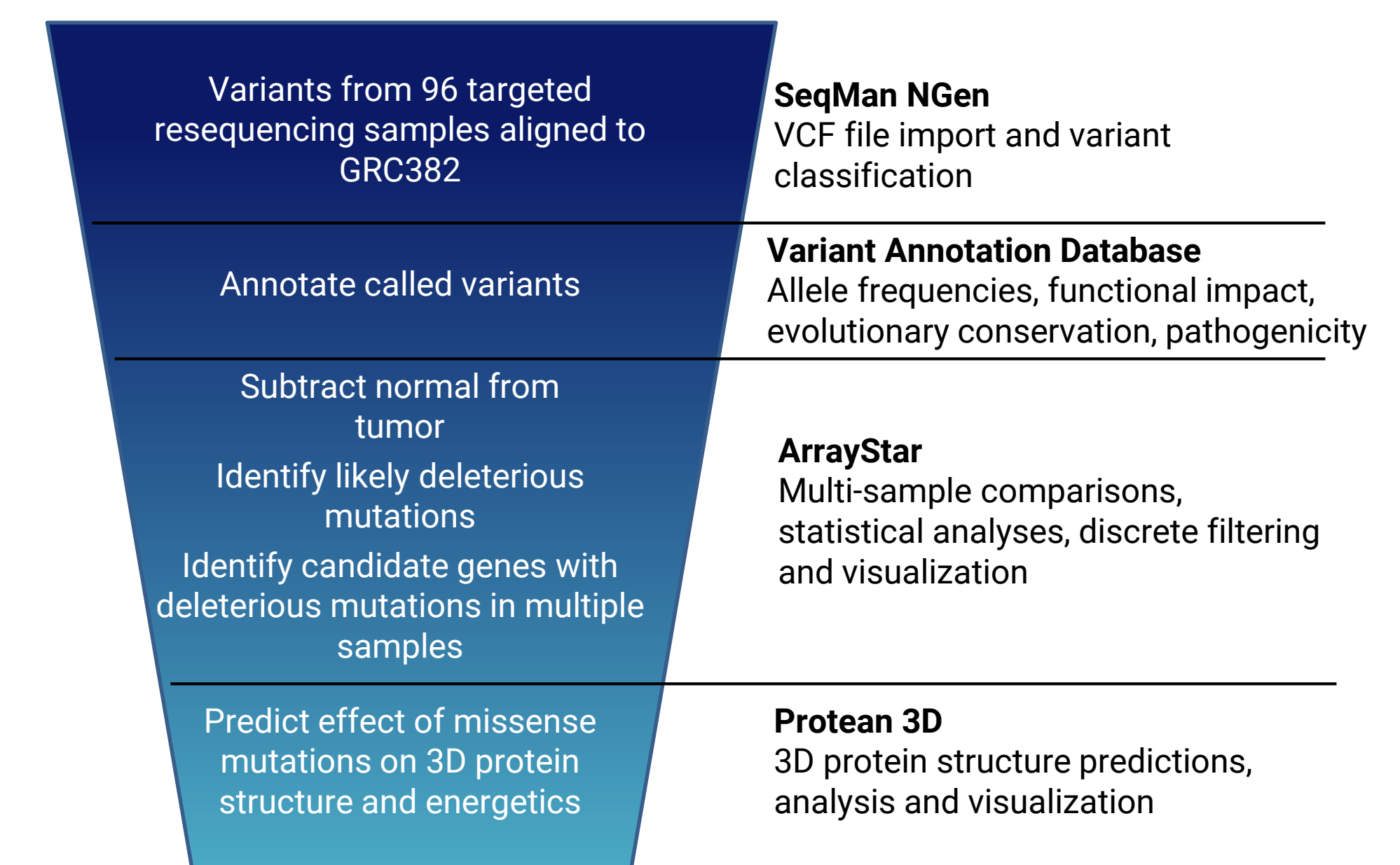
## Project Overview



Figure 1. Overview of DNASTAR Genomics Suite software pipeline for importing VCF variant data, classifying and annotating variants, filtering data, and modeling variants on protein structure, including steps taken to analyze LSCC samples.

## Variant Annotation and Filtering

The DNASTAR Lasergene package can be used to analyze variants from raw sequencing data, or from VCF files produced by other assembly pipelines. Variants are annotated via a two-step process. The first annotation step classifies variants by their effect on coding regions, relative to the imported reference genome. The second annotation step includes import of the DNASTAR Variant Annotation Database (VAD), which combines data from a variety of SNP level annotation databases. The robust downstream filtering options provide researchers efficient means to identify candidate genes based on several criteria including mutation effects, as well as clinical and functional importance.

By applying progressive filters to variant data, we identified a small subset of genes with variants that are predicted to inactivate the gene in multiple tumor samples.
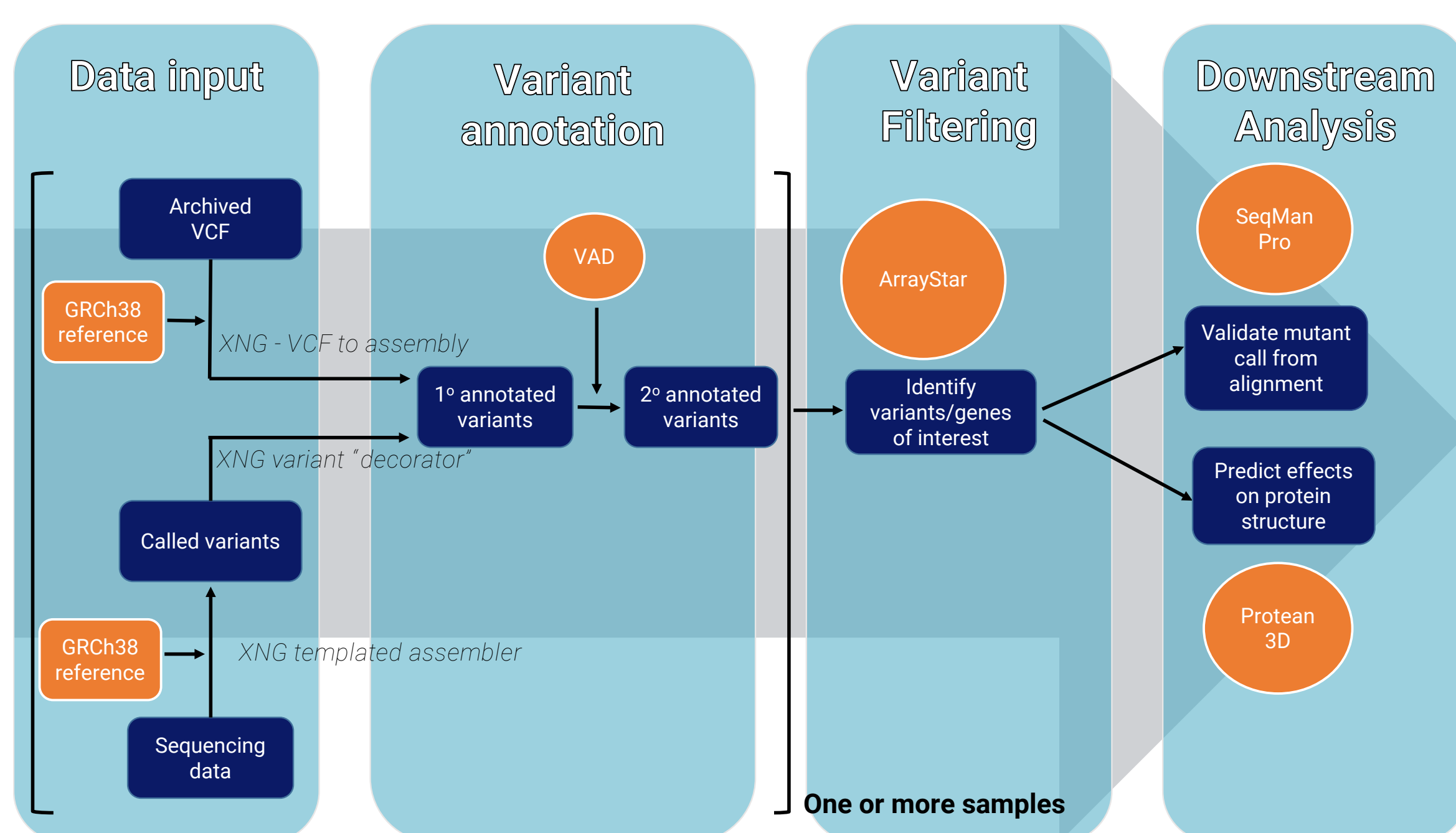


Figure 2. Steps for annotating and analyzing multi-sample variant data in Lasergene. DNASTAR software tools and supplemental data are shown in orange.

**Table 1. Genes with unique nonsense and/or frameshift mutations in multiple samples**

| Minimum number of samples | Number of Genes | Genes |
|---|---|---|
| 1 | 101 | |
| 5 | 6 | CSMD3, FBXW7, LRP1B, PTEN, TP53, TTN |
| 10 | 2 | TP53, TTN |
| 15 | 1 | TP53 |
| 24 | 1 | TP53 |
| 25 | 0 | |

**Table 2. Example of discrete filtering of missense variants in the LSCC samples**

| Step | Criteria | Number of variants | Number of genes |
|---|---|---|---|
| 1 | Missense, minimum variant % = 15, minumum depth = 10 | 3689 | 472 |
| 2 | " + unique to 1 LSCC sample | 2134 | 433 |
| 3 | " + Mutation Taster = Disease causing | 1132 | 303 |
| 4 | " + SIFT = Damaging | 676 | 246 |
| 5 | " + MAF >= 0.01 | 98 | 86 |

**Table 3. Genes with unique deleterious missense mutations in multiple samples[a]**

| Minimum number of samples | Number of Genes[b] | Genes |
|---|---|---|
| 1 | 208 | |
| 5 | 23 | |
| 10 | 7 | CDH10, COL11A1, MDN1, MYH4, NFE2L2, RYR1, TP53 |
| 12 | 2 | MDN1, TP53 |
| 15 | 1 | TP53 |

[a] Starting with 246 gene set from Table 2, Step 4
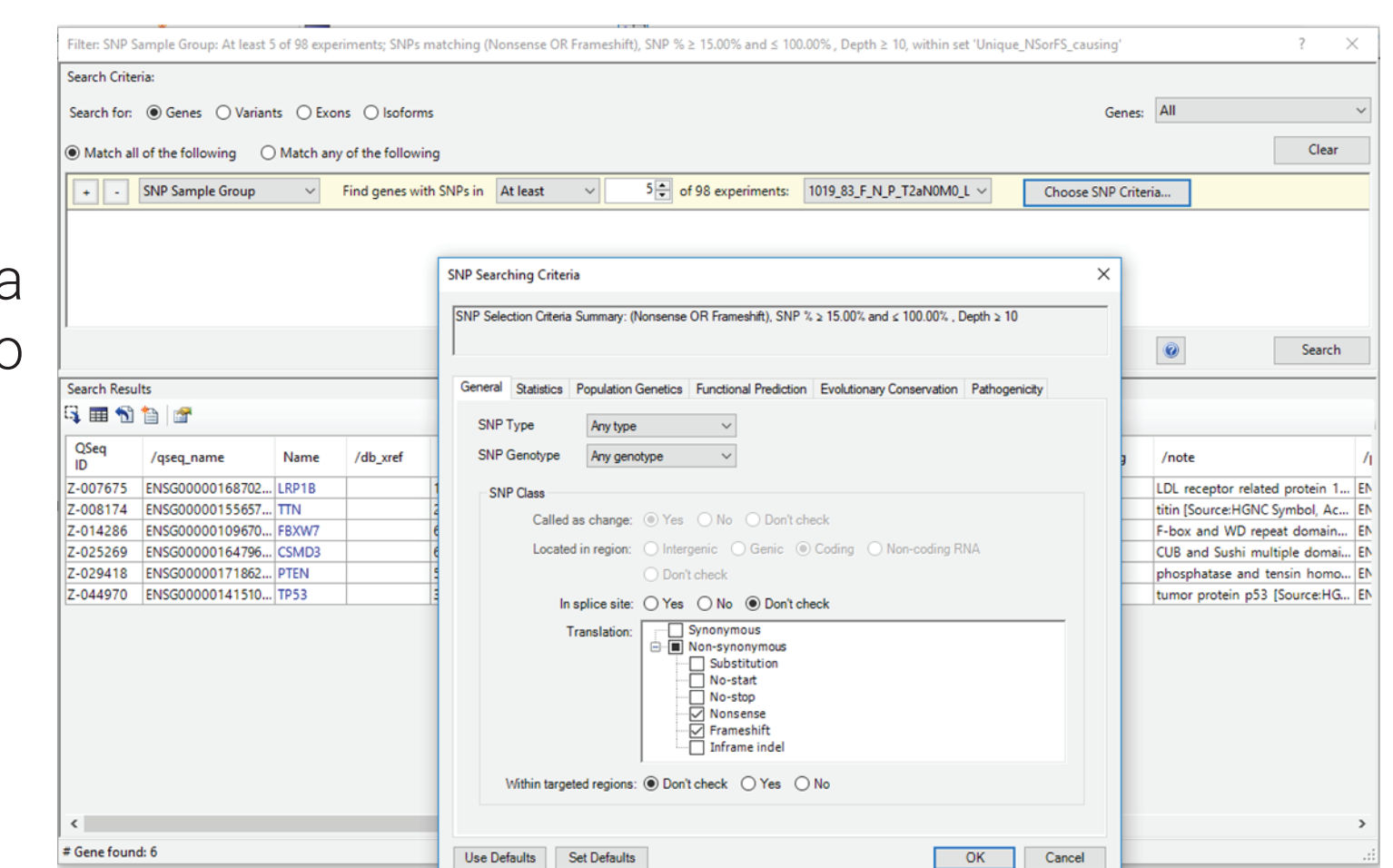[b] Filtered to genes with lengths < 250kb



Figure 3. ArrayStar advanced filtering results showing genes which have unique nonsense and frameshift mutations in multiple tumor samples.
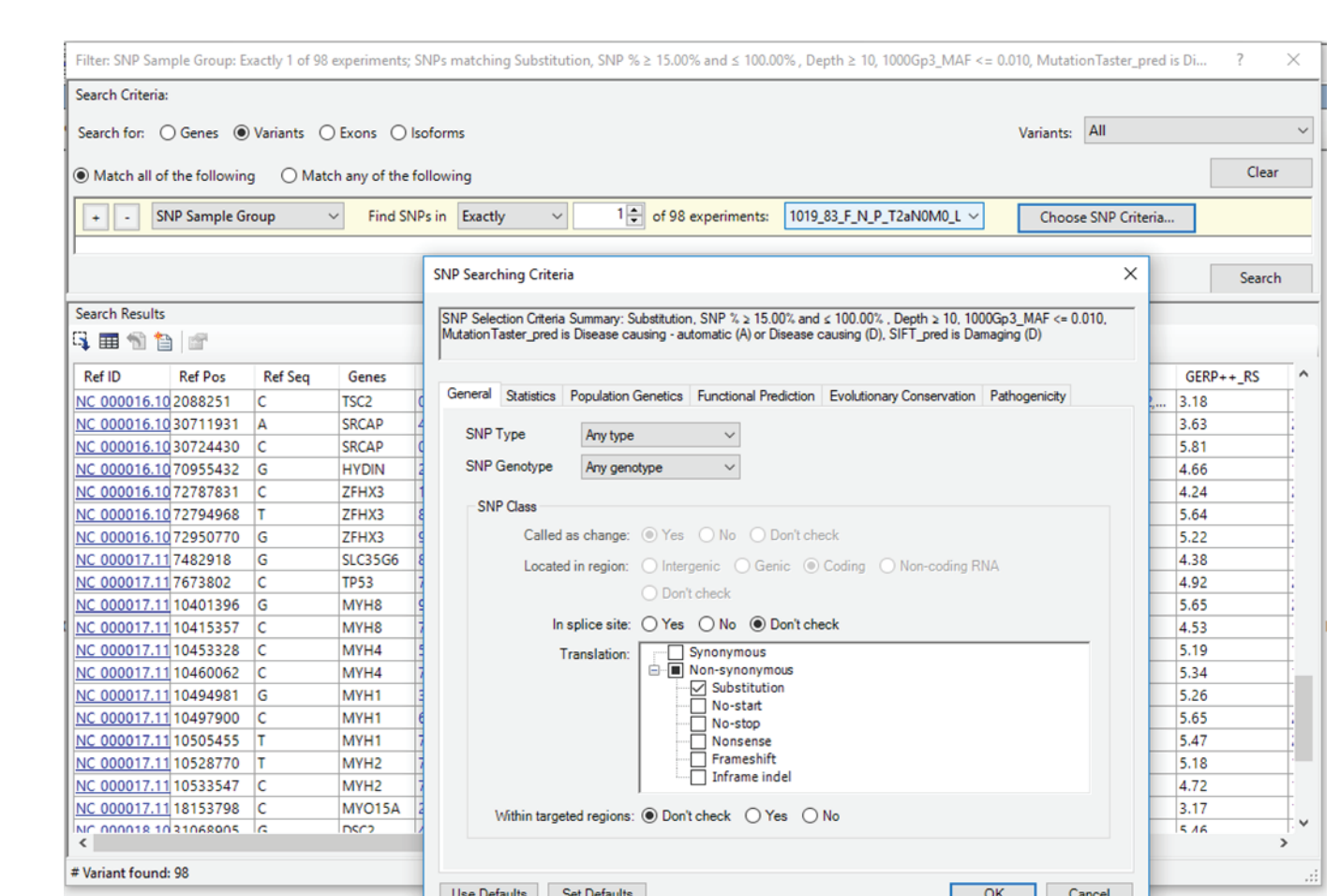


Figure 3. ArrayStar advanced filtering results showing unique missense mutations that are predicted to be deleterious and occur in exactly one sample.

## Structure Visualization and Mutation Modeling

DNASTAR's integrated structure visualization and analysis tools allow users to model mutations of interest on the PDB protein structure. Built-in energy calculations allow users to make guided hypotheses about the effect of mutations on the protein structure and function.

Using the ArrayStar advanced filtering options, we interrogated the TP53 missense mutations to identify those which were predicted to be deleterious by all three functional impact predictors (SIFT, MAF and Mutation Taster) and were predicted to be pathogenic or likely pathogenic in ClinVar. This search yielded three variants, including M2371. Starting with the PDB structure 1TUP[3] which contains TP53 complexed to DNA, we used Protean 3D to mutate the methionine to isoleucine in both the A and B subunits. The structure shows that M237 is in close proximity to the DNA backbone. I237 is predicted to be rotated away from the DNA, a change that may affect binding and would in turn likely be involved in pathogenicity.
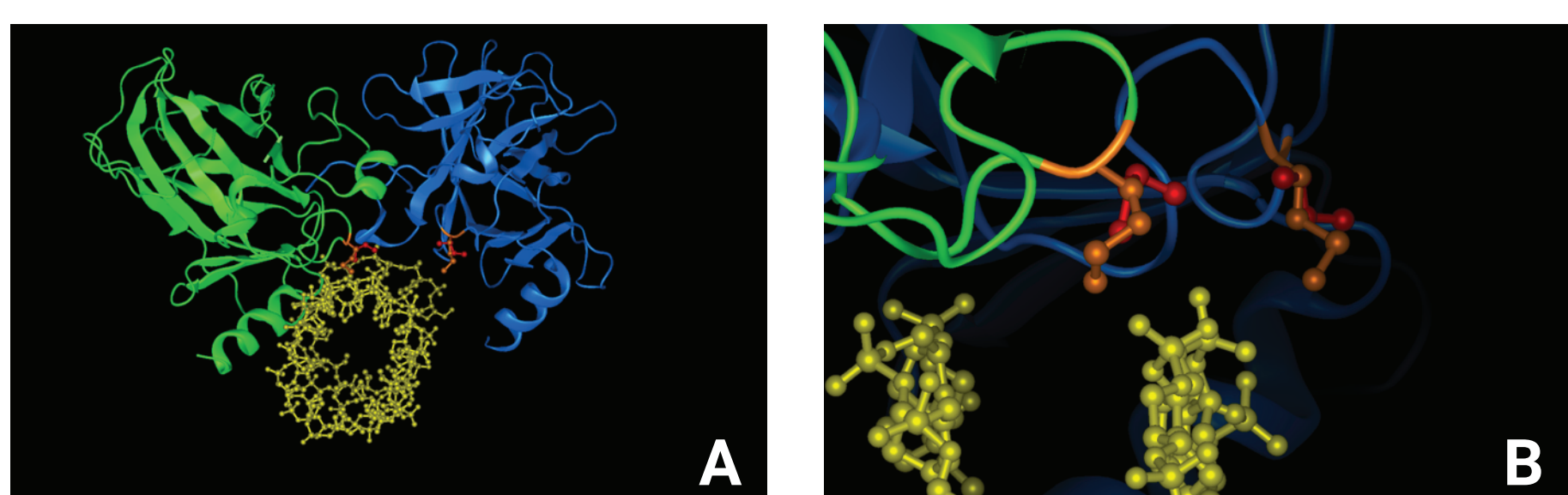


Figure 4. Protean 3D structure prediction of the M237I mutation in TP53 indicates it could affect DNA binding. (A) Crystal structure of a TP53 dimer (green and blue subunits) complexed with DNA (backbone shown in yellow ball and stick rendering). Methionine at position 237 (gold ball and stick) in each subunit was changed to isoleucine and the orientation of the side change (red ball and stick) predicted by the software. (B) Zoomed in image of the methionine (gold) and isoleucine (red) side chains at position 237 of the wild type and mutant, respectively, with the DNA backbone shown in ball and stick rendering (yellow).

## Discussion

The ability to detect and annotate variants from NGS data or VCF files, and integrate the results with multiple variant and gene annotation databases, allows researchers to quickly identify important or interesting mutations for larges number of samples. Fast and intuitive filtering tools allow users to filter on a variety of criteria and combine filtering results in unique ways. Analysis is not limited by the starting data format as long as a suitable reference genome is available.

In addition, by combining structural bioinformatics with sequencing technologies, DNASTAR's integrated workflows can guide genomic and molecular biology researchers to create structure-based hypotheses and to investigate possibilities not evident by sequence data alone.

# www.dnastar.com

References:

1. http://trace.ddbj.nig.ac.jp/DRASearch/study?acc=SRP030634
2. Human genome reference with GRCh38 and ENSEMBL annotations
3. https://www.rcsb.org/structure/1tup