# 7 Steps for Human Variant Analysis

## (and How to Automate Them!)

**DNASTAR**

# ABOUT DNASTAR

DNASTAR is a pioneer in the field of bioinformatics, offering comprehensive software solutions for molecular biology, genomics, transcriptomics, and protein analysis. Our Lasergene Genomics software enables you to set up complex genomic sequencing projects in mere minutes and automates tasks that typically require extensive manual intervention in other software packages.

## About Genomenon

Genomenon is an AI-driven genomics company that organizes the world's genomic knowledge to connect patient DNA to scientific research in the diagnosis and development of treatments for patients with rare genetic diseases and cancer. Genomenon's Mastermind Genomic Search Engine powers variant interpretation by reducing turnaround time, increased diagnostic yield, and lowering the risk of missing a patient diagnosis.

# CONTENTS

# Introduction

In this ebook, we discuss some of the challenges involved in human variant analysis and explore some of the solutions available for addressing these challenges. Whether you are working with a core facility, a bioinformatics group, or doing variant analysis on your own, this ebook will help you learn about important considerations to keep in mind throughout the process.

## Why study variants?

A common workflow in the study of human genetic variation involves the analysis and identification of deleterious variants or of variants associated with a particular population or trait. There are thousands of known variants that cause Mendelian disorders, and thousands more whose molecular basis is yet unknown.

In a typical human variant analysis study, the researcher's goal is to identify which single-nucleotide polymorphisms (SNPs), small insertions and deletions (INDELs), copy number variations (CNVs), or other types of structural variations and rearrangements (SVs) have functional significance. Functionally significant variants are those that cause amino acid changes, abnormal exon splicing, or other protein structure changes that contribute to a diseased state.

## Challenges in variant analysis

The path from raw DNA reads to an understanding of the clinical significance of their variants can be long and complex. Each step and application will affect the accuracy and completeness of the results. An additional complication is that the large number of tools available have led to a myriad of pipelines. Just analyzing one data set can require mastery of up to a dozen bioinformatics applications and online databases.

Another challenge is that many open-source tools are created and then abandoned as graduate students finish their bioinformatics dissertations or laboratories lose funding or change focus. Few of these tools provide comprehensive documentation or instruction.

Even supposedly "automated" open source and commercial pipelines may rely on command line utilities at some of the key steps or require you to navigate through multiple "wizards" at each step in the pipeline. These additional steps generate multiple intermediate data files and increase hands-on time and clock time from raw data to the data analysis steps. It is therefore important to assess your ability and willingness to enter a pipeline that may require the use of command line tools and at least intermediate knowledge of bioinformatics.

# Chapter 1

## Step-by-Step Variant Analysis Procedure

There is no universally agreed upon number of steps involved in variant analysis, and some projects require more steps than others. In this ebook, we have divided the bioinformatics portion of the workflow into **seven steps** based on when you would typically need to switch from one tool to the next to keep advancing through the pipeline.

If you have bioinformatics expertise, you can build these steps into an automated or semi-automated pipeline using scripting and open-source tools. Note that some commercial variant analysis software packages combine two or more of these steps into a single process. If you are working with a core facility, they may have their own pipeline that covers some or all of these steps before delivering the data to you.

If you are new to variant analysis, be sure to do your research and have a plan in place to make sure that the output from one step is readable by the next tool along the path. If a core facility or bioinformatics group is doing some of this work for you, you'll want to understand what data they will be delivering to you so you can plan your next steps accordingly.

# Variant Analysis Overview

**1** Import & clean up sequencing data

**2** Align reads to a reference genome

**3** Remove PCR duplicates

**4** Call variants and INDELS

**5** Filter data to discover important variants

**6** Determine how variants affect genes

**7** Visualize the impact of the variant on the 3D protein structure

# Preliminary Steps

## Choose a sequencing strategy and prepare samples

While this ebook focuses primarily on the bioinformatics and data analysis steps involved in studying human variants, there are many important decisions you will need to make before you can dive into the data. How you design and execute your sequencing strategy will greatly impact the quality of your data, so be sure to consider your analysis goals before beginning your experiment.

Before obtaining DNA samples, first determine if you will be doing a germline or somatic analysis. Somatic analysis compares samples from related tissues from a single individual and is often used in cancer research. Germline analysis, by contrast, compares samples from different individuals that may be heterozygous or homozygous for a trait.

Another consideration is whether you will be doing whole-genome, whole-exome, or RNA sequencing. All of these can yield important results and help elucidate the genetic cause of a disease. RNA-Seq is useful when you want to study differences in gene expression rather than simply the presence or absence of a SNP in the DNA sequence. Whole-genome sequencing (WGS) is not typically used in the case of human subjects. One reason is that it is prohibitively expensive for most researchers, although costs have come down substantially. Another concern is that the significance of non-coding variants is typically much more difficult to determine than coding changes and their abundance (~400x more than in coding regions) makes data handling and analysis more cumbersome and inefficient. Technological improvements and cost reductions are making WGS a more viable option when needed.

Because of the issues mentioned here, whole-exome sequencing, which considers only the ~1-2% of DNA that codes for proteins, is much more common than whole-genome sequencing for variant analysis studies. One caveat to be aware of is that exome capturing kits may miss some regions of interest, potentially missing important variants.

Once you have obtained the appropriate human DNA, each sequencing platform (e.g., Sanger, Illumina, Ion Torrent, PacBio, Oxford Nanopore), will have their own wet-lab protocols for sample preparation specific to the technology. For example, short read sample preparation typically requires fragmentation and PCR amplification steps, while long read preparation avoids DNA fragmentation at all costs. If you would like to enhance mapping accuracy and identify structural rearrangements with more clarity, consider a technology that produces long reads, paired-end, or mate-pair reads.

The final preliminary step is the actual sequencing, which is typically performed at a dedicated sequencing facility.

# Bioinformatics Steps

After sequencing is complete and you have sequence data in FASTQ format or variant call format (VCF), you are ready to start your analysis. The remaining steps all come down to your bioinformatics software and your ability to utilize it. Keep in mind that some of these steps may have been completed before you receive your data.

## Step 1: Import & clean up sequencing data

If your data files are in FASTQ format, this step is where you'll jump in and upload the file(s) into a software application for assessment and possible clean-up.

NGS data files (Illumina and Ion Torrent) are typically cleaned up using the pipeline tools associated with the sequencing instrument. This is normally sufficient. However, some output sequence files can benefit from scanning with a third party tool like FastQC.

**FASTQ files? Here's your starting point**

Sanger data, by contrast, is not typically cleaned up during sequencing. Sanger data usually contains many base calling errors at the 5' and 3' ends where the chromatogram peaks are not high quality. This type of data requires a high-quality software program that can accurately trim the sequence ends,

## Step 2: Align reads to a reference genome

This step involves using a computer program to align/map the read sequences to an existing reference genome; this is usually followed by a local realignment. The two most widely used public sources for human reference genomes are the Genome Reference Consortium (GRC), which provides GRCh36 through GRCh38, and the University of Santa Cruz (UCSC), which provides versions hg18 and hg19.



Some popular tools used for alignment include SOAP, Bowtie/Bowtie2, BWA, and MOSAIK. The results of the alignment are generally saved as BAM or CRAM formatted files.
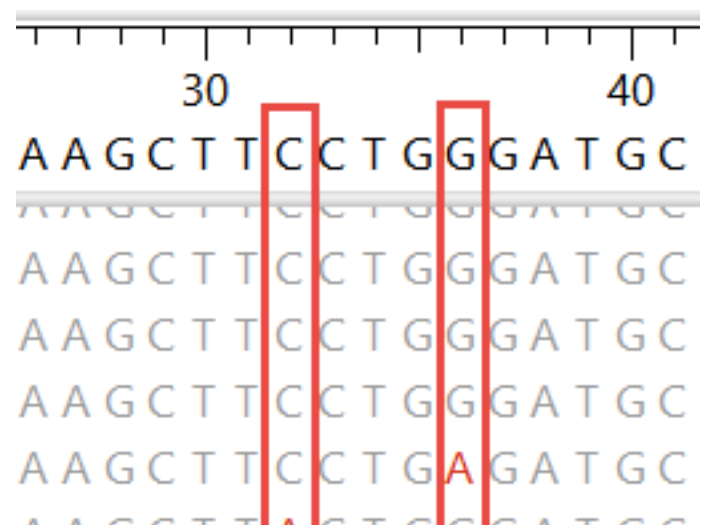
## Step 3: Remove PCR duplicates

When performing whole-genome or whole-exome sequencing, PCR duplicates should be removed immediately after the alignment step. This prevents duplicate reads that originate from a single template from interfering with later variant calling statistics. Commercial solutions may automatically detect and remove duplicate reads from the variant caller. One tool that can be used for this purpose is Picard Tools, though it works only with BAM files.

## Step 4: Call variants and INDELS

Variant calling is the process of comparing the aligned read sequences to the reference sequence in order to find locations where they disagree. The four main types of variants are:



- Single-nucleotide polymorphisms (SNPs)

- Small insertions and deletions (INDELs)

- Larger indels, known as "structural variations" (SVs)

- Copy number variations (CNVs), in which a section of DNA is repeated multiple times

Scores of variant callers are available and are generally classified into four groups:

- **Germline callers** (including CRISP, SAMtools, GATK), typically used for elucidating the causes of rare diseases.

- **Somatic callers** (Including GATK, SomaticSniper), the choice used in most cancer studies.

- **Copy number variation (CNV) finders** (including CONTRA, CNVnator, RDXplorer). CNVs can be detected in both whole-genome and whole-exome assemblies, which is not true of other structural variations.

- **Non-CNV structural variation (SV) finders** (including ExomeCNV, CONTRA) are used to find inversions, translocations or large INDELs.

It is critically important to choose the right variant caller for your data. Using the wrong variant calling pipeline can lead to missed variant calls. In addition, different variant callers have been shown to work better with different types of sequencing technology.

After alignment and variant calling, the limited number of applications that can perform sample cross-normalization or that can utilize a BED file to filter the variant list will perform those steps now.

The output from GATK and other variant calling pipelines is a Variant Calling Format (VCF) file. VCF is a scalable, tab-delimited text file with specific rules pertaining to the order and contents of columns. Each row represents a single nucleotide variant (SNV), insertion, deletion, or other sequence variation. Each variant is uniquely identified with a unique combination of letters and numbers, and the file may contain identifiers from multiple databases.

## Step 5: Filter data to discover important variants

If Steps 1-4 didn't sound familiar or relevant to your experience, that's because they are frequently performed by a core facility or bioinformatics group. This step—which may start with the import of one or more VCF files—is the first one in which most researchers become directly involved.

**VCF files? Jump in here!**

All sequencing platforms produce sequence reads that contain many base-level errors, and it is crucial to apply filtering that can separate "sequencing noise" from the variant "signal." Much of this baseline filtering can be done in an automated fashion, preferably using statistical confidence measures, but some may require you to run some initial "noise filtering" tools.

Even after filtering out sequencing errors, a typical human exome set will still contain thousands of variants. You must now rely on additional criteria, typically imported into the analysis pipeline, to differentiate the uninteresting variants from those that have functional or clinical significance. Some open-source tools that support filtering are VCF Tools and SnpSift.

Many filtering strategies can be used, but here are four examples:

- Filter out **synonymous SNPs.** Since these SNPs encode for the normal amino acid, they will not result in a deleterious effect.

- In the case of a whole-genome study, filter out **variants located in non-coding regions**.

- Filter out **common alleles**. For example, we can reasonably assume that a mutation that causes a rare kidney disease will be rare in the human population. Thus, if a SNP is common (i.e., has a high allele frequency) across a population, it is unlikely to be the cause of the rare disease and can be removed from consideration.

- Some cloud-based tools (e.g., the Integrative Genomics Viewer) and standalone applications (CLC Bio, Lasergene) let you **compare variants from multiple VCF files**. For instance, to find variants related to a particular form of brain cancer, you could filter out variants found in both brain cancer patients and controls and focus on variants found only in the cancer patients.

# Step 6: Determine how variants affect genes

The goal of this step is to determine the functional and clinical consequences of variants and how they impact genes. This is best accomplished with data tables that contain both variant and functional information for one or multiple samples and which have the ability to easily verify variant calls by visualizing the assembled sequence reads. This is most often done using genome browsers, which are available in standalone and web-based versions.

A benefit to **standalone browsers** is that most provide a graphical user interface (GUI) which can support zooming, easier visualization, and the ability to browse interactively. One drawback, however, is that you may be responsible for importing a large number of different data tracks for each of your samples. In addition, cross-comparison and analysis of multiple samples may be challenging, depending on the analysis tools that are available.

One advantage to **web-based genome browsers** (e.g., Ensembl, UCSC Genome Browser, ANNOVAR, AnnTools, VariantAnnotation, NGS-SNP, and snpEff) is that you do not usually need to download or update a human variant annotation database. On the other hand, you must upload your own proprietary data to a remote server, which is not a viable option for all organizations, owing to data privacy policies. Furthermore, some of these tools are command-line driven and can require typing multiple long strings of text into the command-line interface.

At a minimum, a genome browser should allow you to display the aligned reads, see variants clearly (e.g., by displaying them in a different color), view annotation information, and visit the corresponding public database (e.g., Ensembl, the GWAS Catalog, EVA, UniProt, dbSNP, ClinVar) entries online via hyperlinks. Most tools only allow visualization of SNPs, while a few also support viewing CNVs and SVs. Some genome visualization tools allow you to compare sequences from multiple individuals or even multiple organisms.

Once you have tabulated a list of variants, you may identify some for which there is limited data. Most search engines cannot excavate any further than paper titles and abstracts, where 85% of variants are not mentioned. This represents an enormous lost opportunity. While the crowd-sourced ClinVar platform is a useful repository of variant interpretations and related information, it is far from complete. In fact, an average of 31% of submitted references per variant are false-positives, with 30% of submissions having no references at all. Together, the downstream effects of these statistics can be serious, and may result in erroneous conclusions and/or a missed diagnosis. To mitigate these information quality issues, a deeper analysis can be performed using the Mastermind Genomic Search Engine, which is discussed in detail in Chapter 3.

# Step 7 (Optional): Visualize the impact of the variant on the 3D protein structure

An optional step in the variant analysis pipeline is to attempt to determine the effect of a given variant on the 3D structure of the protein it codes for. Sometimes, the 3D structure may be minimally affected, if at all. In other cases, the variant can cause a dramatic change in protein structure, vitally impacting protein function.

Tools like I-Mutant and SDM can predict how the variant might affect protein stability, while structure prediction tools like I-TASSER and Phyre2 can predict the effect on protein structure. However, the latter tools will only work if a structural homolog is available. In some cases, you may find that the structure of the variant protein has already been determined through x-ray crystallography and is available for download through the Protein Data Bank (PDB). You can view these structures through various programs, including the LiteMol viewer.



An easier option is to use DNASTAR's Protean 3D application to accomplish all these tasks. Protean 3D's "Protein Design" workflow lets you use a wizard to easily "mutate" one or more residues and then calculates whether the changes are likely to be stabilizing or destabilizing compared to the original structure. NovaFold, a separately-licensed application that runs through the Protean 3D interface and is powered by the award-winning I-TASSER algorithm, can be used to predict the new protein structure caused by the mutated residue(s). Protean 3D lets you view the original and mutated proteins as fully customizable and rotatable 3D structures.

# Chapter 2

## Using Lasergene Genomics + Mastermind Genomic Search Engine to Streamline Variant Analysis

# HOW IT WORKS

With **Lasergene Genomics** and **Mastermind**, just three applications are required to get meaningful variant results. Once you've identified variants of interest, you have the option to do further downstream analysis and visualization in Lasergene.

## 1 > 2 > 3

**SeqMan NGen**
Import and assemble reads or VCF files with auomatic variant calling.

**ArrayStar**
Filter data to discover important variants and determine how variants affect genes.

**Mastermind**
Determine clinical significance of variants.

> "This is the complete package, from assembly to analysis"
>
> **Marjorie Beggs, Arkana Laboratories**

The **Lasergene Genomics** variant analysis pipeline is completely different from the open-source software or other commercial software described in Chapter 1. How?

Instead of using separate tools and requiring multiple manual interventions at each step, Lasergene fully automates most of these steps into one easy wizard-based application: **SeqMan NGen**. The steps below will take longer to read than to perform, as it typically only takes about **2 minutes** to make selections and begin the assembly and variant calling. After this, variant analysis is performed using **ArrayStar and Mastermind**, with optional analysis and visualizations in other fully-integrated Lasergene applications.

# Getting Started: Assembly and Variant Calling in SeqMan NGen

After launching **SeqMan NGen**, the Workflow screen (**Figure 1**) prompts you to choose a workflow. Select either the Variant Analysis/Resequencing workflow or the RNA-Seq/Transcriptomics workflow if you have raw sequence read data; or the Variant Call Format (VCF) File workflow if your data has already been assembled and the variants detected and saved to VCF format.

**Figure 1.** SeqMan NGen's multi-tabbed Workflow screen is the starting point for the project.

The SeqMan NGen wizard will guide you through the setup for one or multiple samples, auto-grouping and auto-naming FASTQ input files so that assembly (or VCF import), variant detection, variant filtering, sample cross comparison, and variant annotation are all run autonomously. SeqMan NGen supports a wide variety of file types and read technologies.

You can choose to import a custom VCF file and/or BED file to incorporate into your assembly. Adding these files now means the assembled project will include known variants of interest and defined targeted regions, respectively. This information can be used as filtering criteria without requiring you to import additional data after assembly.

Choose to upload a reference sequence or use one of the curated Genome Template Packages provided for common model organisms, including many versions of the human genome (**Figure 2**). The Genome Template Packages include both annotated reference chromosomes and the dbSNP database.

**Figure 2.** A partial list of SeqMan NGen's available genome packages.

| Organism | dbSNP Build | Assembly | Download Size |
|---|---|---|---|
| Apis mellifera | 149 | Amel_4.5 | 129.9 MB |
| Arabidopsis thaliana | 138 | TAIR10 | 87.1 MB |
| Bos taurus | 148 | UMD_3.1.1 | 1.8 GB |
| Caenorhabditis elegans | 138 | WS195 | 54.8 MB |
| Canis lupus familiaris | 146 | 3.1 | 1.1 GB |
| Danio rerio | 142 | Zv9 | 729.0 MB |
| Drosophila melanogaster | 149 | Release_6_plus_ISO... | 108.6 MB |
| Equus caballus | 151 | EquCab2.0 | 1.1 GB |
| Escherichia coli K12 MG1655 | - | ASM584v2 | 3.1 MB |
| Gallus gallus | 147 | 5.0 | 625.6 MB |
| Homo sapiens | 150 | GRCh37.p13 | 4.9 GB |
| Homo sapiens | 150 | GRCh38.p7 | 3.5 GB |
| Homo sapiens - Ensembl w/PDB | 150 | GRCh38.p10 + Ense... | 3.5 GB |
| Homo sapiens mitochondrion | 150 | GRCh38.p7 | 33.3 kB |
| Homo sapiens, CEU reference | 142 | GRCh37.p10 | 3.8 GB |
| Homo sapiens, CHBJPT reference | 142 | GRCh37.p10 | 3.8 GB |
| Homo sapiens, YRI reference | 142 | GRCh37.p10 | 3.8 GB |

Select    Cancel

In the Analysis Options screen (**Figure 3**), choose variant calling options (diploid, haploid, somatic), gender, and initial variant filtering. Check the variant annotation database (VAD) box to automatically annotate any found variants with data from multiple databases, including Mastermind citation counts and hyperlinks to the Mastermind Genomic Search Engine.

**Figure 3.** SeqMan NGen "Analysis Options" wizard screen.



You can also opt to detect CNVs and/or SVs. Choosing these options produces an assembly project that includes assembled data, filtered variant calls, CNVs, SVs, and extensive variant annotations.

Follow the guided steps to initiate assembly locally or on the cloud. Once assembly is complete, proceed to open the project in **ArrayStar** to analyze the results.

**Figure 4.** ArrayStar's "Experiment List" screen.

# Find important variants with ArrayStar

**ArrayStar automates Steps 5-6** with versatile filtering tools and rich graphical and tabular displays showing the variant results in context. Start by using the Advanced Filtering tool to filter out uninteresting variants, leaving only those worthy of further study. See the next chapter for an example of how to apply variant filtering.

Next, move to the Variant Table (**Figure 5**), which is ready for analysis and allows you to effortlessly add (or remove) data columns with information about each variant. If you checked the VAD box in SeqMan NGen, there are hundreds of columns to choose from. Many of these include links to the database entry for the variant in databases like the Mastermind Genomic Search Engine, ClinVar, dbSNP, 1000 Genomes Project, the Gene Ontology Consortium (GO), GERP, MutationTaster, PhastCons, dbNSFP, and many others.

ArrayStar is the only tool we are aware of that additionally performs an automatic cross sample analysis. If a variant is detected in one sample, the application cross checks all other samples at the same position and reports the status (variant, no variant call, no coverage) as a Variant Table column

**Figure 5.** ArrayStar's Variant Table displaying a selection of customizable data columns.

# Further exploration of variants with Mastermind

Also accessible from ArrayStar is **Genomenon's powerful Mastermind Genomic Search Engine,** a comprehensive genomics database designed to be an essential resource for analysts faced with data uncertainty. Mastermind allows you to quickly and easily search and cross-reference variant data from millions of PubMed articles. Additionally, Mastermind provides entries for millions of variants that are not included in databases such as HGMD and differentiates them with much higher precision.

Mastermind enables access to millions of full-text articles encompassing the entire breadth and depth of published genomic evidence. Using patented algorithms, Mastermind allows you to filter and prioritize relevant articles with high specificity, without compromising sensitivity. Additionally, Mastermind is able to draw genetic associations between keywords (e.g., genes, variants, diseases, CNVs, phenotypes, and therapies) and deliver a segmented literature output that has been customized for your particular case.

The successful use of this tool should result in a highly useful yield of evidence surrounding your variant of interest. On the other hand, if your results are limited in Mastermind, you can be assured that limited data exists.

# Optional Analysis in Lasergene

After identifying important variants using ArrayStar and Mastermind, you can use other integrated Lasergene applications to visualize your assembly read alignment and coverage and/or model mutations on protein structure (see Step 7 in Chapter 2).

To conclude this section, we have provided a table contrasting the time and number of steps to perform variant analysis with DNASTAR software and with competing software. In the next chapter, we'll show you an example of how to combine the filtering in ArrayStar and Mastermind to home in on variants of interest in clinical sequencing data.

**Table 1.** Typical steps involved in the variant analysis workflow using DNASTAR versus other software tools, with average elapsed time to complete the workflow, including both hands-on and compute time.

| Other Pipelines | Lasergene Genomics and Mastermind |
|---|---|
| **1. Import & clean up sequencing data**<br>• Import FASTQ<br>• Assess data & remove flawed portions | **Assembly and Variant Detection in SeqMan NGen** |
| **2. Align reads to a reference genome**<br>• Align reads to a reference genome<br>• Local realignment | |
| **3. Remove PCR duplicates** | |
| **4. Call variants and INDELs**<br>• Call variants<br>• Call INDELS<br>• Remove marginal variants<br>• Import VCF and identify known variants<br>• Import allele and genotype frequencies<br>• Predict amino acid changes<br>• Import variant databases, including Mastermind Genomic Search Engine<br>• Multiple functional predictions | |
| **5. Filter data to discover important variants** | **Variant Analysis in ArrayStar and Mastermind** |
| **6. Determine how variants affect genes**<br>• Cross-sample analysis<br>• Determine clinical significance using Mastermind and other databases | |
| **7. Visualize the impact of the variant on the 3D protein structure (optional)** | **Mutation Modeling in Protean 3D** |
| **Total average time: >4 hours**<br>**10+ applications** | **Total average time: 2 hours**<br>**3 integrated applications** |

# Chapter 3

## Use Case with Lasergene Genomics & Mastermind Genomic Search Engine

A recent collaboration between DNASTAR and Resonant Therapeutics used RNA-Seq Illumina reads to study gene expression differences in ovarian carcinoma cells (OVCAR-2 cell line) grown using 2D and 3D cultures. In case these terms are unfamiliar, a 2D culture is where cells are grown in a petri dish, while a 3D culture is grown on a matrix that mimics the microenvironment of a tumor. In this study, we show how to quickly reduce a large data set down to three variants of interest, and then demonstrate the utility of Genomenon's Mastermind Genomic Search Engine for exploring deleterious variants.

# Setting up and running the assembly in SeqMan NGen

SeqMan NGen was used to align Illumina RNA-Seq data from multiple replicates to the GRCh38. p2 human genome package. In addition to read alignment, the following steps were performed automatically during the assembly process:

1. SNPs and other small variants were called using the diploid model

2. Variants were annotated using the Variant Annotation Database (VAD)

3. Differential gene expression values were calculated using DESeq2

This study involved cancer cells, which understandably have a large number of mutations. During assembly, SeqMan NGen called 965,209 variants.

# Filtering Variants in ArrayStar

After assembly in SeqMan NGen, we used the Advanced Filtering feature in ArrayStar to identify variants found in all six experiments that met the following criteria (**Figure 6, next page**):

• General tab - Non-synonymous SNPs only

• Statistics tab - Pnotref > 90%, SNP% > 15, Depth > 20

• Pathogenicity tab - Labeled in the ClinVar database as Pathogenic or Likely Pathogenic

**Figure 6.** Selections made in ArrayStar's "Advanced Filtering" dialog.



**Figure 6a**. Filtering dialog



**Figure 6b**. General filtering options with non-synonymous filter applied.



**Figure 6c**. Statistical SNP filtering options.



**Figure 6d**. Pathogenecity filtering options with ClinVar filter applied.

ArrayStar found **three variants** that fit all criteria **(Figure 7).**

**Figure 7.** The three variants of interest found after filtering in ArrayStar.



Looking at these results in the Variant Table, we were able to see the SNP base and reference bases, the number of articles found in Mastermind, links to that evidence, the clinvar_trait, and the MutationTaster_pred. For legibility in this ebook, the Variant Table has been split in half to create two figures (**Figures 8 and 9**).

**Figure 8.** The leftmost columns in ArrayStar's Variant Table.



**Figure 9.** The rightmost columns in ArrayStar's Variant Table, along with a tooltip for the selected variant and a view of the Details panel.

# Exploring Genomic Associations with Mastermind

Note that the variant in the KRAS gene (middle of the three variants) has over 4,300 Mastermind citations. This well-known variant, located on chromosome 12 at position 25245350, causes the G12V change in the KRAS protein and is the key driver mutation behind the tumors being studied. In addition to the two-star status of this SNP in ClinVar, notice that all three functional prediction methods present in the VAD also flag it as deleterious. As shown near the bottom of the tooltip, this KRAS gene mutation has been associated with a wide variety of cancers.

ArrayStar includes a direct link to Mastermind, where we can view the entry for known KRAS mutations (**Figure 10**).

**Figure 10.** The Mastermind Professional Edition landing page for the KRAS G12V variant.

# Optimize output with prioritized sorting and filtering

The image on the previous page shows the full list of over 4,300 publications citing this variant. The literature base can then be focused by adding phenotype keywords. In this case, when "ovarian carcinoma" is applied, the reference number is reduced to 140 articles (**Figure 11**).

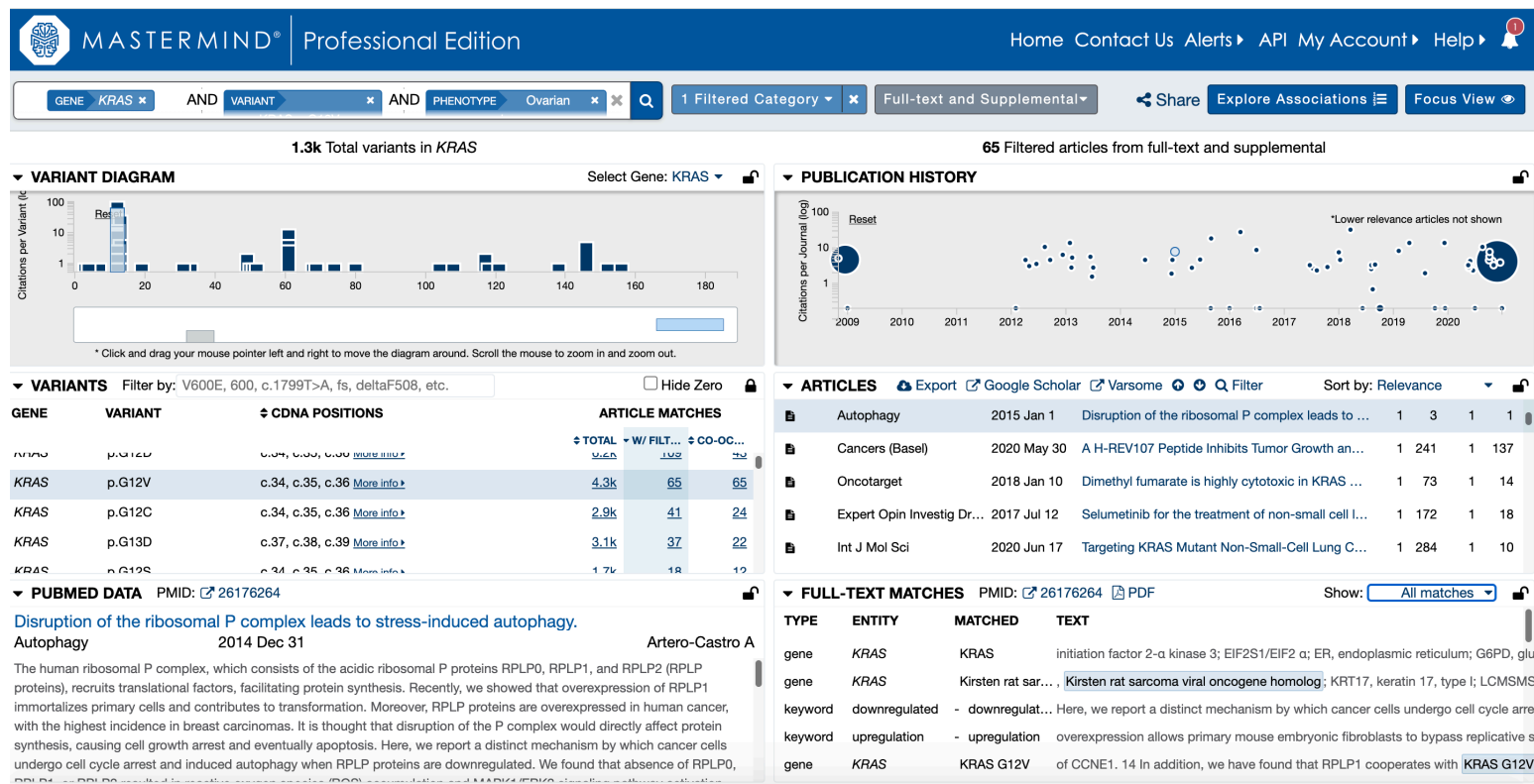**Figure 11.** Mastermind KRAS G12V with ovarian carcinoma phenotype applied.



Other filters may be used to further refine this search, including expression filters in the "Genetic Mechanism" tab (**Figure 12**).

**Figure 12.** Mastermind Filters including expression

This focuses the results on **65 highly relevant articles** (**Figure 13**).

**Figure 13:** Mastermind KRAS G12V Filtered by phenotype, expression and relevance



The articles displayed are sorted by relevance, with the larger circles in the publication history pane denoting the most relevant articles. Relevance is determined by several factors, including, but not limited to:

- The number of times in which the gene ("KRAS") and variant ("G12V") names are mentioned in each article,
- Whether these terms are in the title, abstract and/or full text or supplemental information, and
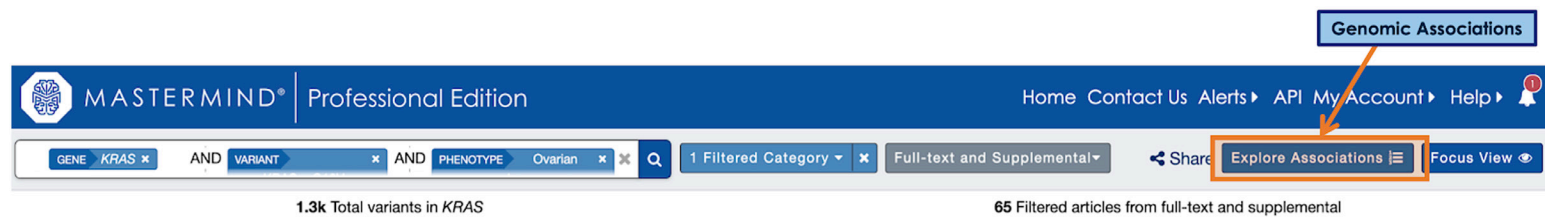- Proximity to phenotype and other keywords applied through the filtering process.

All of these elements can be viewed in the "Full Text Matches" section by viewing all matches being highlighted in the sentence fragments.

## Uncover data connections with genomic associations

Due to the highly complex nature of genetic disorders and the variety of factors that contribute to their pathogenicity, the opportunity to create clarity around how your search criteria of interest relate with each other is exceedingly beneficial. The Genomic Associations functionality within Mastermind allows you to explore significant genetic associations that may otherwise be missed or unrecognized, giving you the ability to uncover key data that may lead to a new hypothesis, a research discovery, or even solving a clinical case.

To expand on the example, suppose that you are interested in exploring treatment options for ovarian carcinoma related to KRAS p.G12V mutations. Clicking on the blue "Explore Associations" box (**Figure 14**) directs you to a flexible interface that connects your search keywords with all associated genetic factors in the medical evidence.
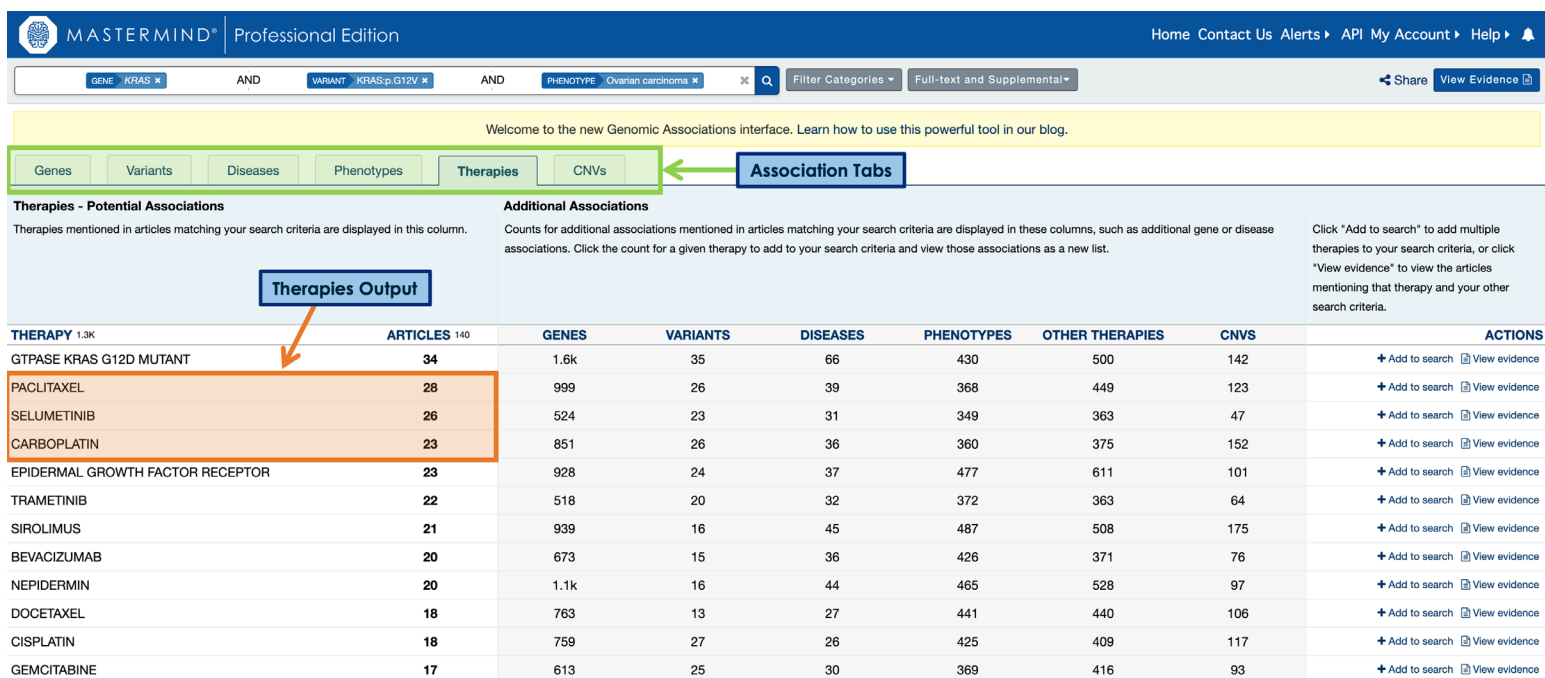
**Figure 14.** The Explore Associations box on the Mastermind page.



In this case, the "Therapies" tab reveals prioritized articles covering therapeutic intervention for KRAS p.G12V-related ovarian carcinoma. As seen in Figure 15, in terms of article quantity, the top three pharmacotherapies are:

1)    Paclitaxel
2)    Selumetinib
3)    Carboplatin

**Figure 15.** The "Therapies" tab with the top three pharmacotherapies boxed in orange.



Awareness of therapies surrounding a particular genetic variant creates informed focus during hypothesis formation and study design. Most importantly, however, this invaluable insight significantly impacts clinical decision making by adding precision to patient care.

Finally, although these patient-facing applications are compelling, they represent a small fraction of this technology's complete potential: as evidenced by the other five association tabs, there's much more for Lasergene and Mastermind users to discover.

# Conclusion

Next-generation sequencing (NGS) of whole genomes, exomes, and RNA-Seq data is a powerful tool in biomedical research and diagnostics. However, sifting through the vast amount of NGS data to find the most relevant variants can be like looking for a polar bear in a snowstorm. Incorrect or deficient annotations can cause you to overlook potentially disease-causing variants. Open source and commercial tools have been developed to address this issue, but often, a dozen or more unrelated tools must be mastered in order to perform the analysis.

In response, DNASTAR has developed Lasergene Genomics, a fully integrated variant analysis and annotation pipeline with an intuitive, easy-to-use interface. Lasergene Genomics and its Variant Annotation Database greatly simplify variant analysis so that you can get accurate answers without needing a PhD in bioinformatics. With the recent addition of the powerful Mastermind Genomic Search Engine, the once time-consuming job of variant analysis is now faster, easier, and more accurate than ever before.

# CONTACT US

To see a demo of the software, or to learn more about how Lasergene Genomics and Mastermind can fit into your research, email or call us at one of the numbers listed below.

www.dnastar.com

Email: sales@dnastar.com

Phone: 608.258.7420

Toll Free in the U.S. and Canada: 866.511.5090

Call free from the U.K.: 0.808.234.1643

Call free from Germany: 0.800.182.4747