# DNA-Seq Analysis

## Introduction

This tutorial describes the DNA-Seq analysis workflow in Avadis NGS using an amplicon sequencing dataset. It assumes that Avadis NGS is installed and the steps listed in the 'Getting Started' tutorial have been completed.

## Dataset

This is a paired-end amplicon sequencing data generated from 96 Yoruba HapMap samples. For the purpose of this tutorial, we will consider three samples from the same family (father, mother and child) out of the original 96. The dataset zip file (`dnaseq-small-dataset-illumina.zip`) contains the BAM files from these three samples. The reads were aligned against the hg19 assembly.

## Goals

In this tutorial, we will check the quality of the amplicon resequencing visually and quantitatively, perform SNP detection and visualize the results, and priortize the list of identified SNPs for further analysis.

## Sample Import

For this tutorial, we need to create a DNA Variant Analysis experiment with hg19 assembly and UCSC annotations. Choose "Illumina" as the sequencing platform and "Paired End" as the library layout.
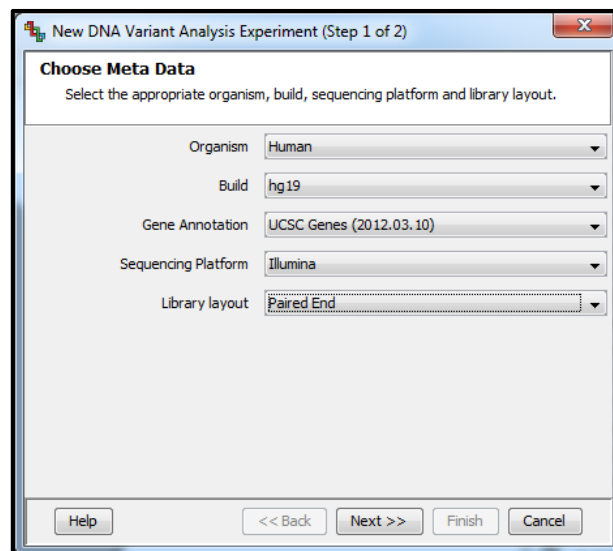


**Figure 1: DNA-Seq experiment creation**

While loading the BAM files present in the zipped file `dnaseq-small-dataset-illumina.zip` (using the 1st icon from the left), we have an option of renaming the samples to make them easily distinguishable.
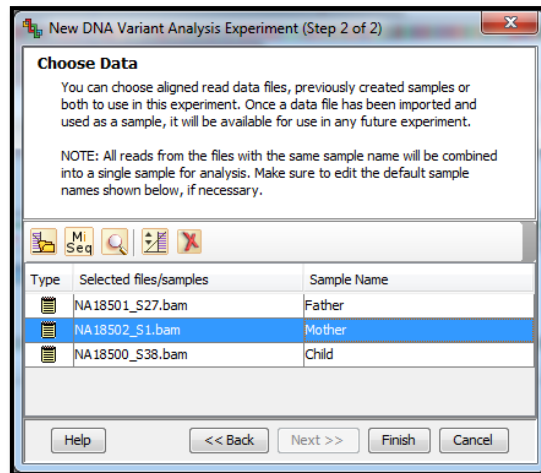
Strand Life Sciences
Algorithms for Life

**Figure 2: Renaming samples**

Once the samples get loaded, the genome browser is launched showing the chr1 of each sample. The fact that similar peaks are present in each sample reflects the fact that this is indeed amplicon data.



**Figure 3: Genome browser tracks on initial launch**

## Quality Control

Since we are planning on running SNP analysis, the quality control metrics need to be carefully considered using the quality inspection in the workflow navigator (using the Quality Inspection section in the workflow browser).

## Pre-Alignment QC Plots

Using these plots, we can discover the base quality by its position in the read, the approximate GC% of the amplicon regions, and the average quality of reads and bases in each sample.
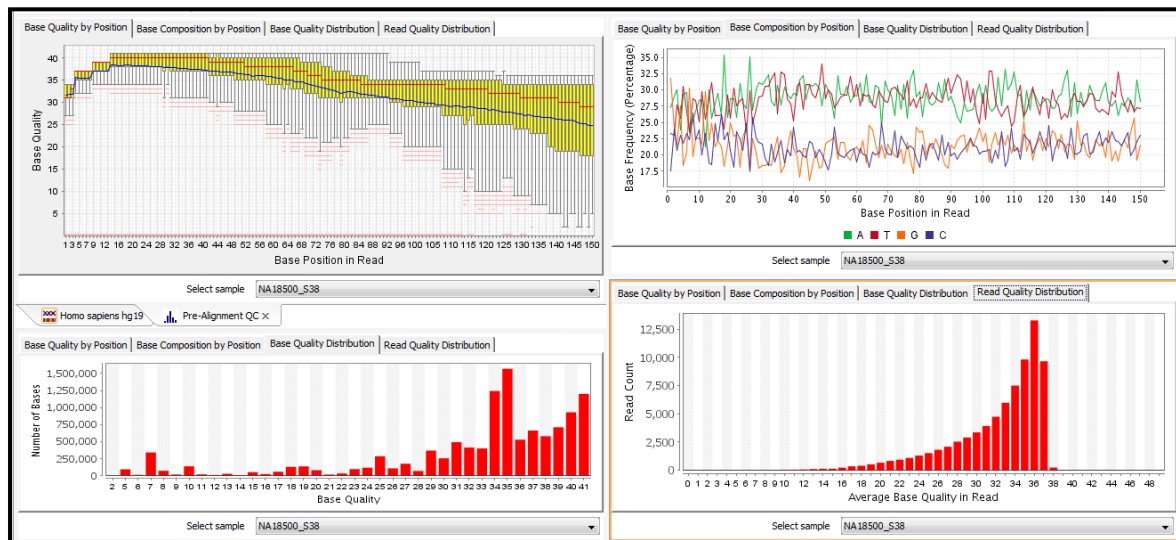


Figure 4: Pre-Alignment QC plots

We can also take a look at QC plots specific to Illumina data (using Quality Inspection → Base Quality by Tile workflow step). The plot shows average base qualities rendererd as a heatmap with user controls for choosing the sample and lane of interest. Rows in the heatmap correspond to tiles of the lane under consideration, while columns represent flow cycles. Each cell is colored according to the average quality of all the bases that fall in the specific cycle of the specific tile.

The below figure shows the Base Quality QC plot for this dataset in which the qualities are good in the earlier flow cycles and degrade marginally towards the end.
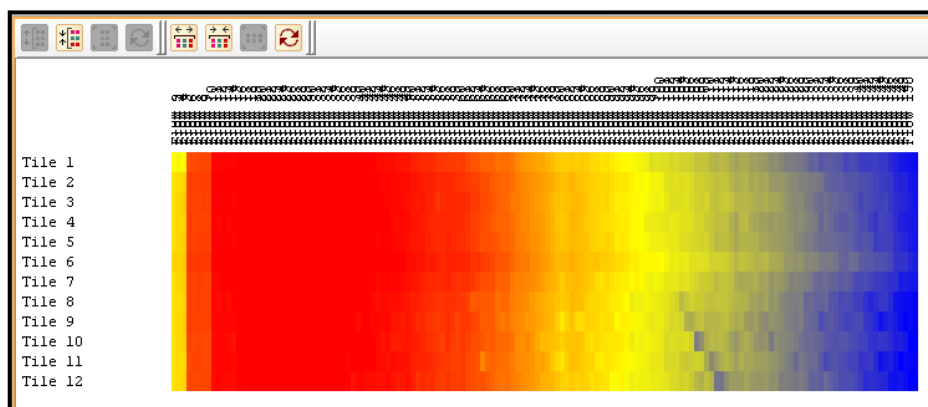


Figure 5: Base quality by tile

## Post Alignment QC

A paired end library of a normal sample (i.e. not a cancer sample) is considered good quality if the majority of pairs align in the expected orientation and distance from each other. We can check the 'Match status' pie-charts (using the Quality Inspection → Alignment QC Plots workflow step), to verify if the paired end libraries are good.
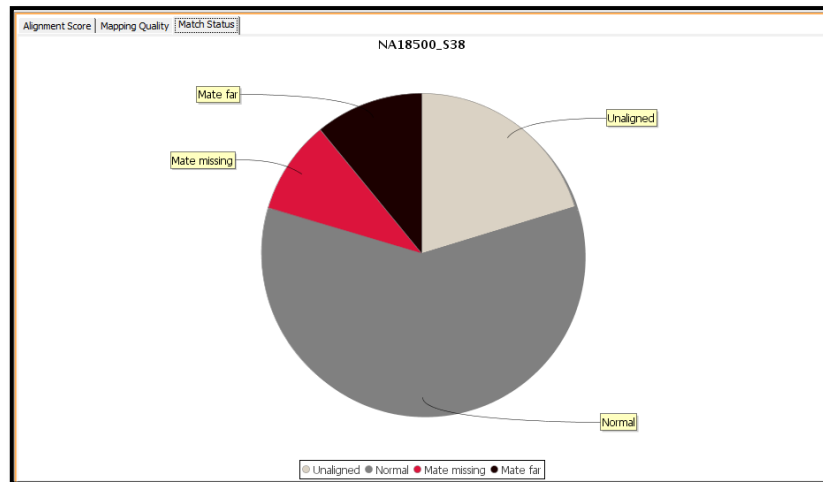
**Figure 6: Match status of paired end reads**

## Targeted Region QC

One of the most important QC steps in amplicon analysis is to determine the efficacy of targeted resequencing. For this, we need to import the file containing information about the target region first (using the Utilities→ Import Region List workflow step). The Region list is a simple list of chromosome, start, end, strand and other annotation information, such as zygocity, score, type etc. For this tutorial, use the `Target-Regions.tsv` file provided with the dataset. (Note that currently, this is obtained by extracting the information from the [Targets] section of the Manifest file present in the MiSeq folder). After giving the file location in the first step, we can proceed with the defaults till step 4 and then mark the mandatory columns as shown below.



**Figure 7: Column selection**

After the region list gets imported, we can use this as input for the Quality Inspection → Targeted Region QC workflow step and run this step with default parameters.
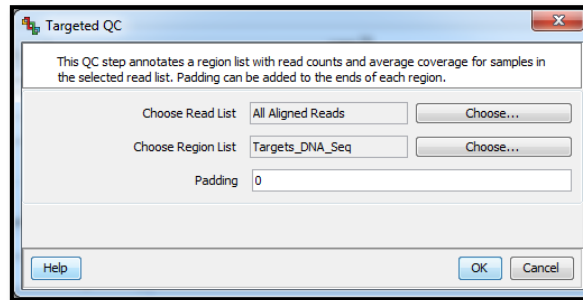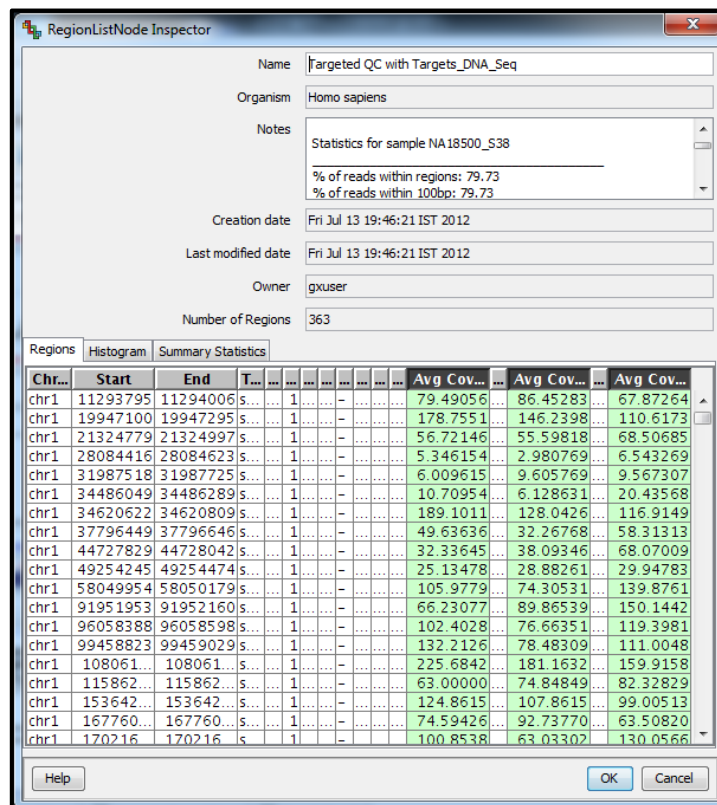


**Figure 8: Target region QC**

For this tutorial dataset, the coverage across the target regions is very high with all three samples having coverage of above 75%. This can be determined by going through the 'Notes' section in the inspector of the output of this analysis step.



Another important QC step is to determine if there are any regions with uniformly low coverage across all the samples. To do this visually, we can perform clustering on both samples and regions based on the average coverage (using Utilites→ Cluster Regions workflow step). This would give us the below image wherein we can see certain regions having uniform low coverage (blue) across the samples.
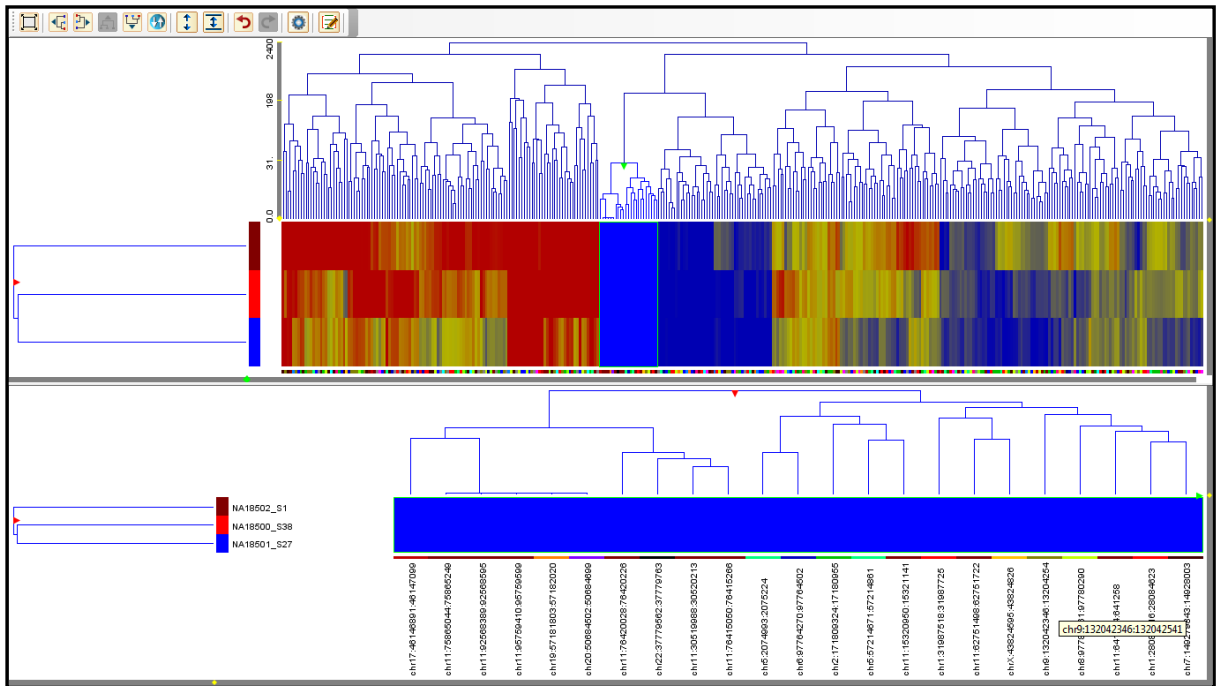
Figure 9: Cluster on target region QC list

## Filter

Before proceeding with SNP Detection, we can filter out some of the low quality reads (using Filter → Filter by Read Metrics workflow step) with criteria as shown below. This would leave us with a total of 156,512 reads from the original total of 229,262 reads.
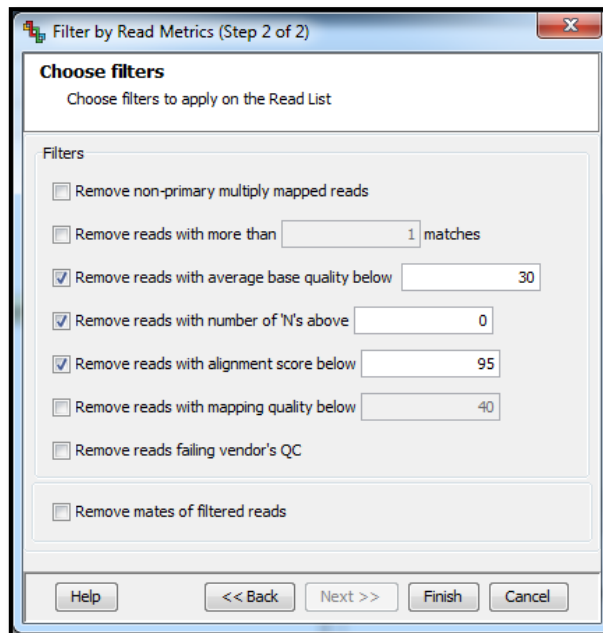


Figure 10: Filter criteria

## SNP Detection

After QC and filtering, we are ready to proceed with the variant analysis. We can run SNP Detection on the filtered read list using the Analysis→ SNP Detection workflow step. We can choose a dbSNP annotation database. This needs to be downloaded from the 'Annotations Manager'. Please note that this takes a significant amount of time so we can choose the option as 'None' if we want to skip that part.



**Figure 11: SNP Detection Parameters**



**Figure 12: SNP results**

The SNP results are divided into single-base variant lists and multi-base variant lists. The SNPs of each sample are listed independently. In addition a multi-sample report is created which is the union of all variants found across all the samples. We can double-click and open each of these outputs to get an idea of their contents.
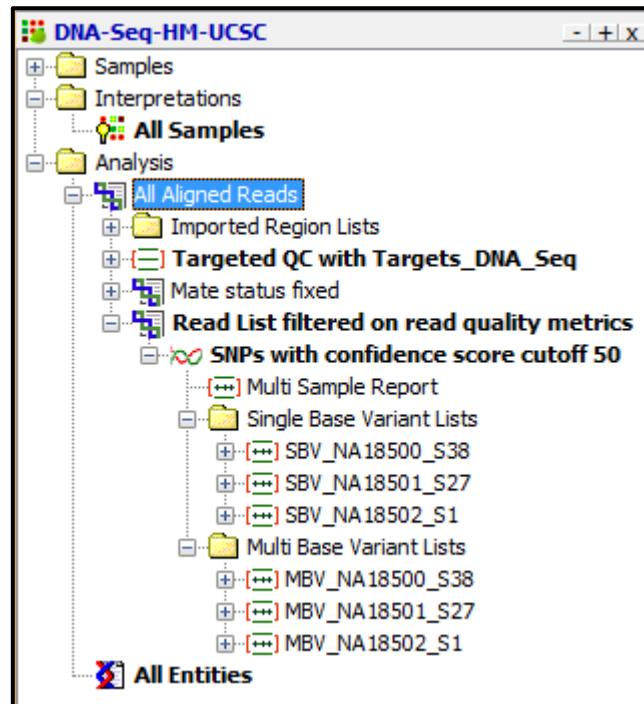


**Figure 13: Experiment navigator after SNP calling**

The multi-sample report has one variant allele per row. For each variant allele it has six sample specific columns:

- total reads
- % reads supporting variant allele
- % reads which are different from the reference
- Strand bias
- SNP call at this location (could be ref)
- Score of the SNP call (empty if no SNP)

Since we have 3 samples, we will have 3x6 plus the additional allele specific columns in the multi-sample report. The next step would be to narrow down this list of SNPs to the most interesting ones according to the use case.

## Find Significant SNPs

In this step, we need to specify the multiple sample report and an interpretation depending on the experimental setup. The below figure shows the various experimental setups that are supported in this analysis step.
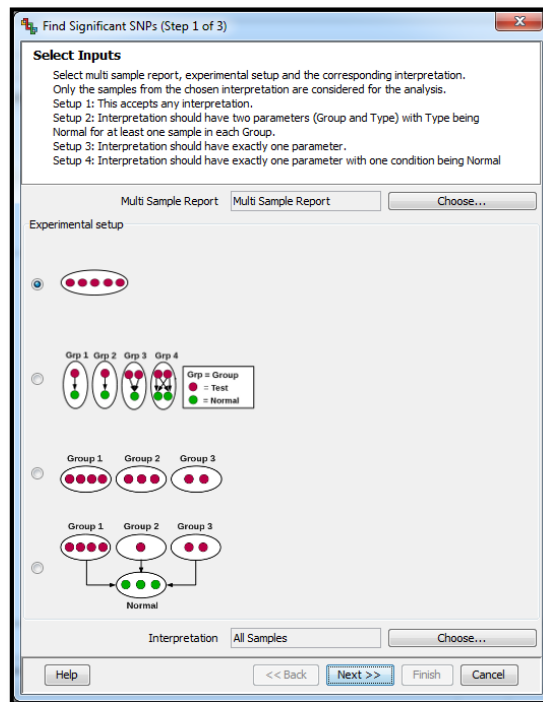
Figure 14: Find Significant SNPs

On selecting the first setup, 119 of the 346 alleles listed in the multi-sample report are present in 3 out of 3 samples with reasonable number of supporting reads (35%). The number of supporting samples can be changed using the slider.
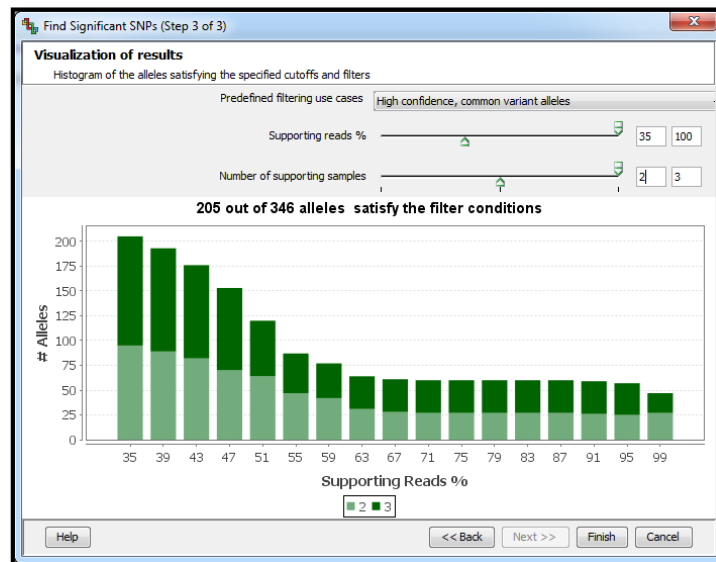


Figure 15: Visualization of results

## Visualization of Results

To quickly browse through the 119 alleles that were saved, we can right-click on output in the experiment navigator and select the option to navigate in the Genome Browser. The Navigate Spreadsheet tab in the Genome Browser will show the contents of the region list.



**Figure 16: Navigating the region list in the Genome Brower**

When we double click on a specific region as shown in the below image, the view zooms to that location. We can also look for a particular region of interest using the search functionality. Alternatively we can also scan through the browser looking for interesting SNPs. In the below figure, for this particular variant, the parents are heterozygous whereas the child is homozygous indicating a loss of heterozygosity.
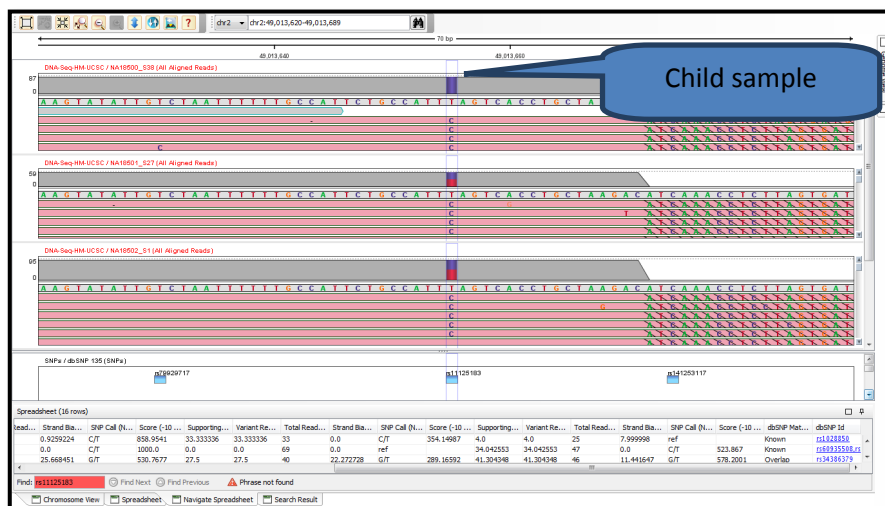


**Figure 17: GB view of the SNP**

From the mismatch histogram alone it might be hard to glean information about the alleles present at this location.  For this, we can right-click on the genome browser track at the location of interest and launch the variant support view. It confirms that in the present case it is a homozygous variant in the off-spring.
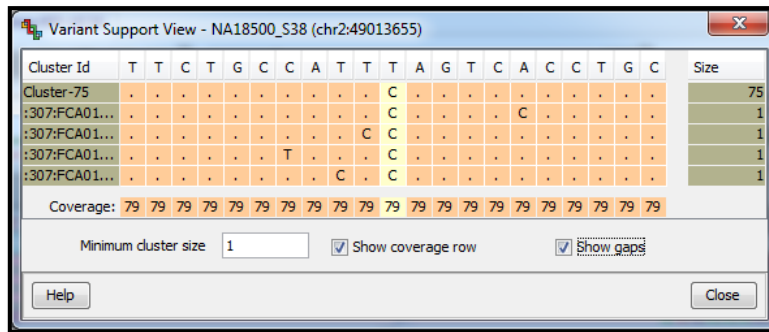
**Figure 18: Variant Support View**

## SNP Effect Analysis

We can see the effect of these 119 alleles by using the SNP effect analysis. Please select all the effects and run the analysis. This would output a region list as well as a gene list. In the present case, we would see the presence of a SNP in an exonic region and the resultant affected gene(MPZL1) in the gene list . This functionality is similar to the Ensembl SNP Effect predictor tool.
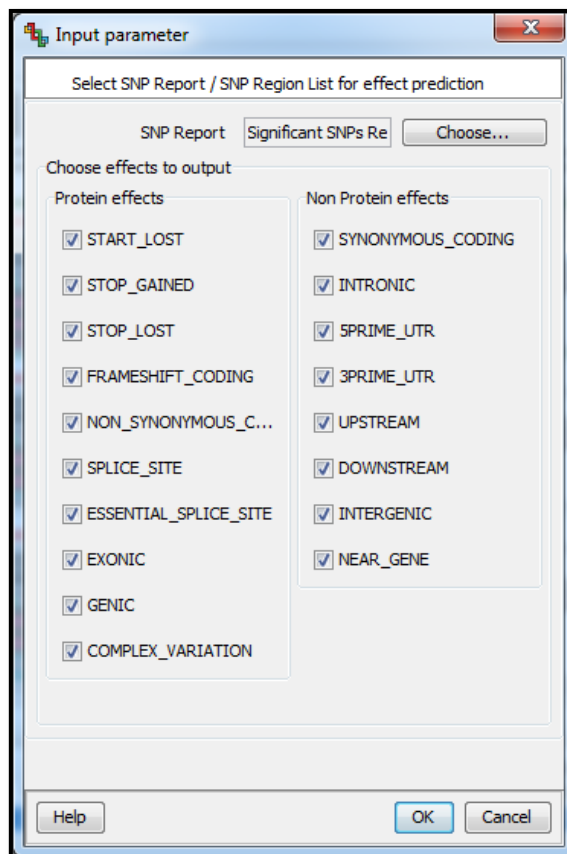


**Figure 19: SNP effect analysis**

If you have access to the whole dataset containing the MiSeq run folder for all the 96 samples, you can load that in and repeat this analysis with the larger dataset.

This concludes our DNA-Seq tutorial. This is a very brief overview of the DNA-Seq experiment workflow in Avadis NGS. For more details or clarifications, please revert back (sales@avadisngs.com or support@avadisngs.com) and we will address your queries.