

Galaxy / vSNP

Greg consulted with Dr. Vivek Kapur in the Huck Institutes of the Life Sciences to build the Galaxy / vSNP analysis environment. Whole genome sequencing for disease tracing and outbreak investigations is routinely required for high-consequence disease. The Galaxy / vSNP tools are critical components of several Galaxy workflows that locate and validate SNPs in bacterial samples and produce annotated SNP tables and corresponding phylogenetic trees. The tools are capable of processing very large sets of inputs and efficiently accommodate multiple genome references for mapping samples.

The Galaxy / vSNP toolset consists of the following tools which are components of larger Galaxy workflows.

- **Determine reference from data** - sniffs bacterial genome samples to discover the primary species within the sample
- **Add zero coverage** - adds zero coverage to a VCF file to prepare it for processing by the Get SNPs tool
- **Statistics** - produces an Excel spreadsheet containing statistics for samples and associated metrics files
- **Get SNPs** - collects quality parsimonious SNPs from VCF files and outputs alignment files in FASTA format
- **Build tables** - produces annotated SNP tables in the form of Excel spreadsheets from outputs produced by the Get SNPs tool

The Galaxy / vSNP environment provides some very interesting tools and workflows. Bacterial samples often consist of a “mixed bag” of species, so it's not always clear to the researcher which is the prominent species in the sample. This makes selecting the correct reference genome for mapping difficult or impossible.

The Galaxy vSNP: determine reference from data tool inspects the bacterial sample data and looks for what we call “dnprints” that match with the same dnprints in one of many reference datasets installed into the Galaxy / vSNP environment. The reference dataset containing the same dnprints will be selected for mapping. The tool outputs a dataset consisting of only the Galaxy dbkey string (metrics and other information is available in a separate output), and this dataset is used as input to one of Galaxy's workflow parameter tools which can be used to specify the selection of the dynamic select list that the Galaxy mapping tools use for reference genome selection.

Figure 7 shows the workflow that performs this DNA sniffing, reference selection and mapping process.

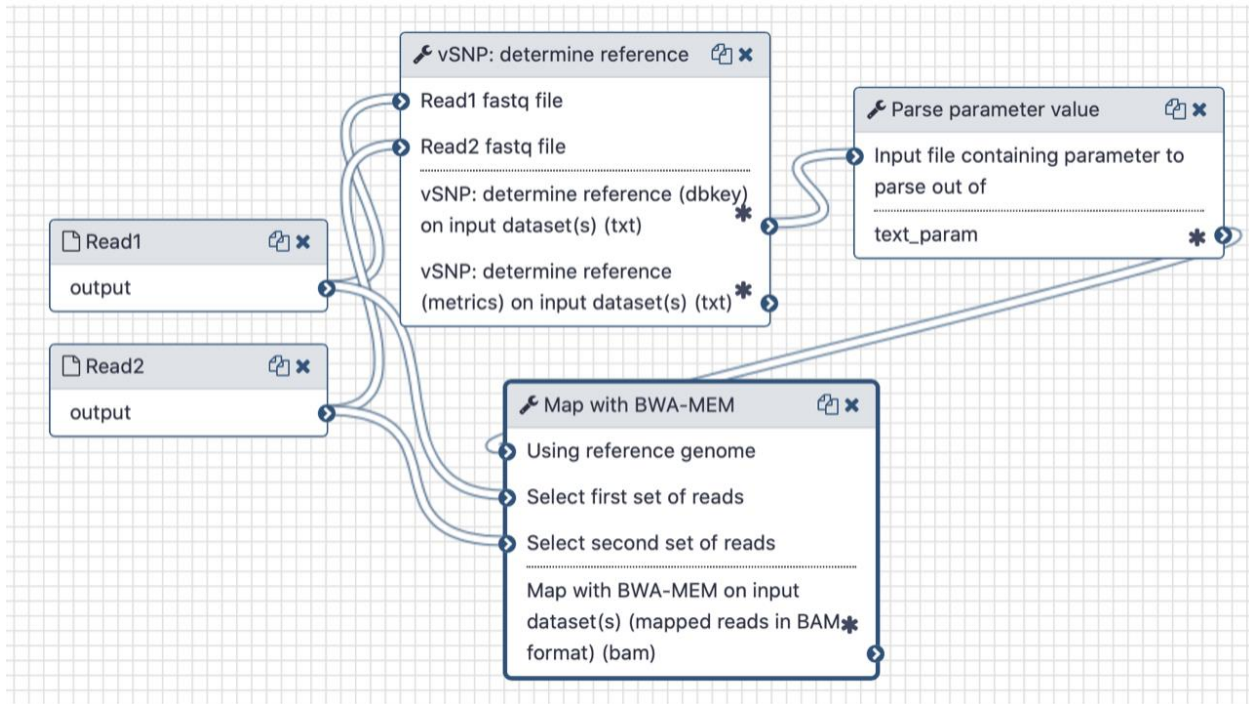


Figure 7: vSNP: determine reference from data, the first DNA sniffer tool for Galaxy

Figure 8 displays an example of an annotated SNP table produced by the vSNP: build tables tool. The tool accepts a combination of single SNPs JSON, average MQ JSON and NEWICK files (or associated collections of each) to produce these annotated SNP tables. The SNPs JSON and average MQ JSON files are typically produced by the Galaxy vSNP: get SNPs tool and the NEWICK files are typically produced by the Galaxy Phylogenetic reconstruction with RaXML tool.

The SNP tables display closely related isolates and enable identification of mixed SNPs when multiple bacterial strains are infecting an organism. The table structure is shown in Figure 8. The columns identify the genome location of the SNP calls and the isolates are contained within the rows. The reference (or ancestral strain if the reference is an outgroup) is listed across the top, identified as the "reference call". SNPs that are not highlighted will match the reference. The map-quality row values are the average of the map quality scores of each isolate in that position. These scores measure the confidence that the read has been mapped to the correct

location on the genome. The maximum score possible is 60, and lower scores lessen the confidence that the SNP was correctly identified. The annotation of the position is provided at the bottom of the table.

The image shows a table of SNP data with several columns highlighted by red boxes and arrows. The columns are: 'reference_pos', 'reference_call', 'L3-1057_MI_Alco_Beef_27', 'L3-3082_MI_Midl_Beef_59-J_1', 'L3-1950_MI_Saga_Dairy_56-A_9', 'L4-2093_MI_Saga_Cat_56-A', 'map-quality', and 'annotations'. The 'reference_pos' column contains NCBI reference numbers like 'NC_002945.4-1799442'. The 'reference_call' column contains nucleotide calls like 'C'. The 'L3-1057...' row contains nucleotide calls for different samples, with some cells highlighted in green or red. The 'map-quality' column contains the number '60'. The 'annotations' column contains text descriptions of the genomic regions, such as 'possible transcriptional regulatory protein, None, BQ2027_MB0202'.

reference_pos	reference_call	L3-1057_MI_Alco_Beef_27	L3-3082_MI_Midl_Beef_59-J_1	L3-1950_MI_Saga_Dairy_56-A_9	L4-2093_MI_Saga_Cat_56-A	map-quality	annotations
NC_002945.4-1799442	C	T	T	T	T	60	POSSIBLE CONSERVED MEMBRANE PROTEIN, None, BQ2027_MB1636
NC_002945.4-232188	G	N	C	C	C	60	possible transcriptional regulatory protein, None, BQ2027_MB0202
NC_002945.4-3413355	G	A	A	A	G	60	possible triacylglycerol synthase (diacylglycerol acyltransferase), None, BQ2027_MB0202
NC_002945.4-3464485	C	T	T	T	T	60	PROBABLE PYRUVATE FORMATE LYASE ACTIVATING PROTEIN PFLA (FORMATE)
NC_002945.4-546200	G	N	G	A	A	60	PROBABLE CTDOCHROME P450 137 CYP137, BQ2027_MB3710C
NC_002945.4-546200	T	C	T	T	T	60	PHOSPHOTYROSINE PROTEIN PHOSPHATASE PTPA (PROTEIN-TYROSINE-PHOSP
NC_002945.4-3685739	G	T	T	T	T	60	hypothetical protein, None, BQ2027_MB3370
147961-4-546200_CN	A	A	G	G	G	60	POSSIBLE INTEGRAL MEMBRANE PROTEIN, None, BQ2027_MB1778C
546812-4-546200_CN	T	A	A	A	A	60	PROBABLE TRANSCRIPTIONAL REGULATORY PROTEIN, None, BQ2027_MB1966
060642-4-546200_CN	C	C	T	T	T	60	HYPOTHETICAL METHYLTRANSFERASE (METHYLASE), None, BQ2027_MB0214C
5696406-4-546200_CN	G	G	A	A	A	60	HYPOTHETICAL PROTEIN, None, BQ2027_MB2796C
8288848-4-546200_CN	C	C	T	T	T	60	No annotated product
1482792-4-546200_CN	C	C	T	T	T	60	HYPOTHETICAL PROTEIN, None, BQ2027_MB3359C
211749-4-6-6-172	G	G	A	A	A	60	No annotated product
NC_002945.4-2523406	G	G	A	A	A	60	No annotated product
4281207-4-2071827	A	A	A	A	A	60	Probable glycine dehydrogenase gcvB (Glycine decarboxylase) (Glycine cleavage
NC_002945.4-2332060	G	A	A	T	T	60	transmembrane serine/threonine-protein kinase jpkri (protein kinase j) (stpk j)
NC_002945.4-2279531	G	G	A	A	A	60	No annotated product

Figure 8: vSNP: build tables produces annotated SNP tables in the form of Excel spreadsheets

Figure 9 displays an example of a phylogenetic tree that would be associated with a SNP table similar3 to the one shown in Figure 8.

Phylogenetic Tree from Best-scoring ML Tree:

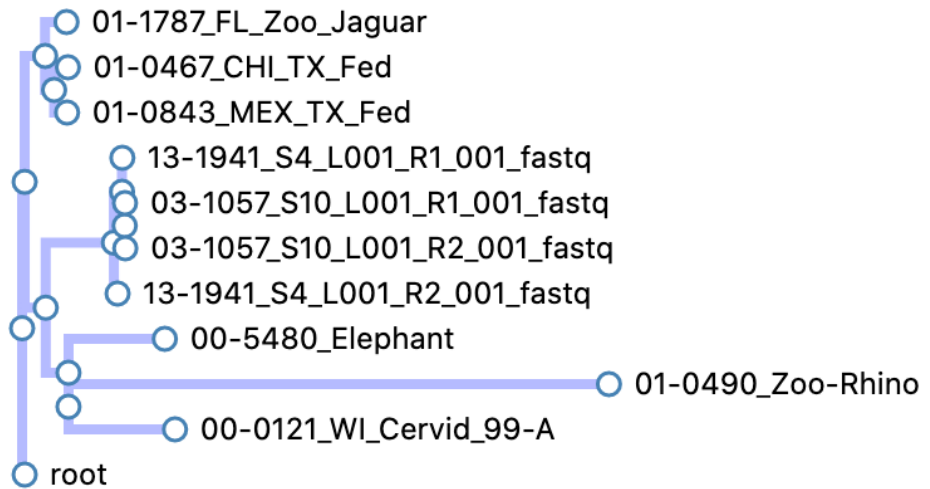


Figure 9: Phylogenetic tree produced by the Galaxy Phylogenetic reconstruction with RaXML tool