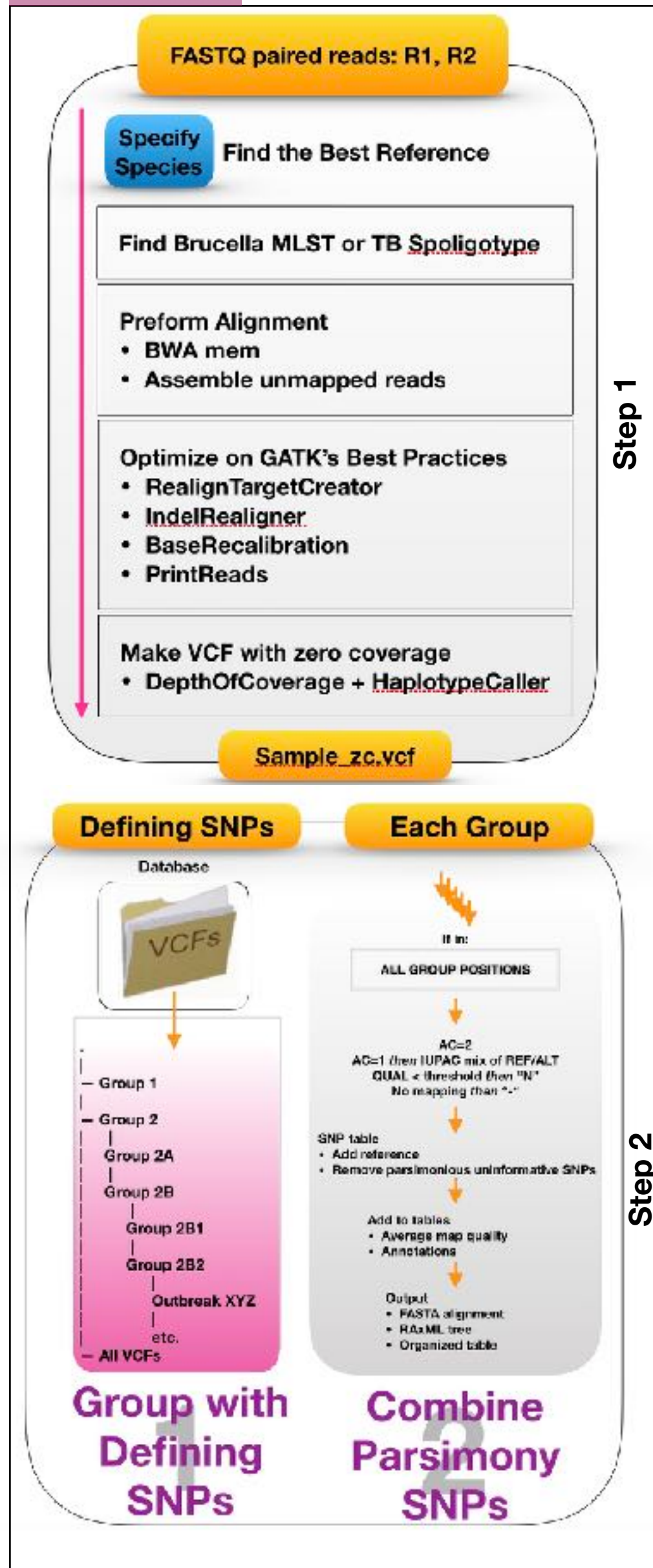# vSNP: Bacterial validation SNP genotyping tool utilizing Mycobacterium and Brucella species

Tod Stuber[a], Suelee Robbe-Austerman[a]

[a]United States Department of Agriculture (USDA), Animal & Plant Health Services (APHIS), National Veterinary Services Laboratories (NVSL), Diagnostic Bacteriology Laboratory, Ames Iowa, USA. [b]USDA, APHIS, NVSL, Diagnostic Virology Laboratory, Ames Iowa, USA.
Corresponding authors (515) 337-7388. (tod.p.stuber@usda.gov, suelee.robbe-austerman@aphis.usda.gov )

**Barcode**

## PIPELINE



FASTQ paired reads: R1, R2

Specify Species — Find the Best Reference

Find Brucella MLST or TB Spoligotype

Preform Alignment
• BWA mem
• Assemble unmapped reads

Optimize on GATK's Best Practices
• RealignTargetCreator
• IndelRealigner
• BaseRecalibration
• PrintReads

Make VCF with zero coverage
• DepthOfCoverage + HaplotypeCaller

Sample_zc.vcf

Step 1

Defining SNPs | Each Group

Database

VCFs

Group 1
Group 2
Group 2A
Group 2B
Group 2B1
Group 2B2
Outbreak XYZ
etc.
All VCFs

ALL GROUP POSITIONS

AC=2
AC=1 then IUPAC mix of REF/ALT
QUAL < threshold then "N"
No mapping than "-"

SNP table
• Add reference
• Remove parsimonious uninformative SNPs

Add to tables
• Average map quality
• Annotations

Output
• FASTA alignment
• RAxML tree
• Organized table

Step 2

Group with Defining SNPs **1** | Combine Parsimony SNPs **2**

## Overview

Single nucleotide polymorphism (SNP) analysis of high-throughput sequencing (HTS) data is increasingly the preferred method to genotype bacterial outbreaks. A specific tool was needed to rapidly call, validate and compare SNPs from FASTQ files in a timely manner while utilizing large datasets. vSNP was developed to address these challenges. It is particular well suited for outbreaks or clonal organisms such as Mycobacterium tuberculosis and Brucella species.

vSNP is publicly available at: https://usda-vs.github.io/snp_analysis/

## Method

vSNP runs on macOS and Linux systems. It is written in Python 3, utilizing Python packages and system programs available from Anaconda package manager. The program is run from the command-line, requiring minimal experience to setup and execute.

vSNP provides a high-resolution SNP analysis from Illumina pair FASTQ files. vSNP is implemented in a two-step process. Step 1 chooses the best reference, and outputs alignments and SNP calls into VCF files. Step 2 takes in a collection of VCF files and outputs SNP alignments and phylogenetic trees. Predetermined SNP positions are used to create groups. Grouping samples on defining SNPs provides manageable datasets and allows results to be quickly and thoroughly validated. Poor SNP positions are easily targeted for further visual validation using Excel formatted SNP alignment tables. Phylogenetic trees are created from validated tables. SNP tables provide the average map quality and annotation at each position.

## Conclusion

To achieve high resolution genotyping, it is necessary to validate SNPs within an outbreak. This must be accomplished quickly to minimize personnel time and to not delay actionable results. We have found it best to validate using SNP tables with fewer than 50 SNP positions. Defining SNPs are used to group similar samples into smaller, more manageable groups. This use of a defining SNP to group closely related isolates is a key element that differentiates vSNP from other pipelines. It allows both a perspective against other isolates and allows the validation necessary to make confident, actionable decisions.

Other clonal organisms, or bacteria or even viruses within an outbreak situation, can also be incorporated into vSNP to provide the same high-level resolution.
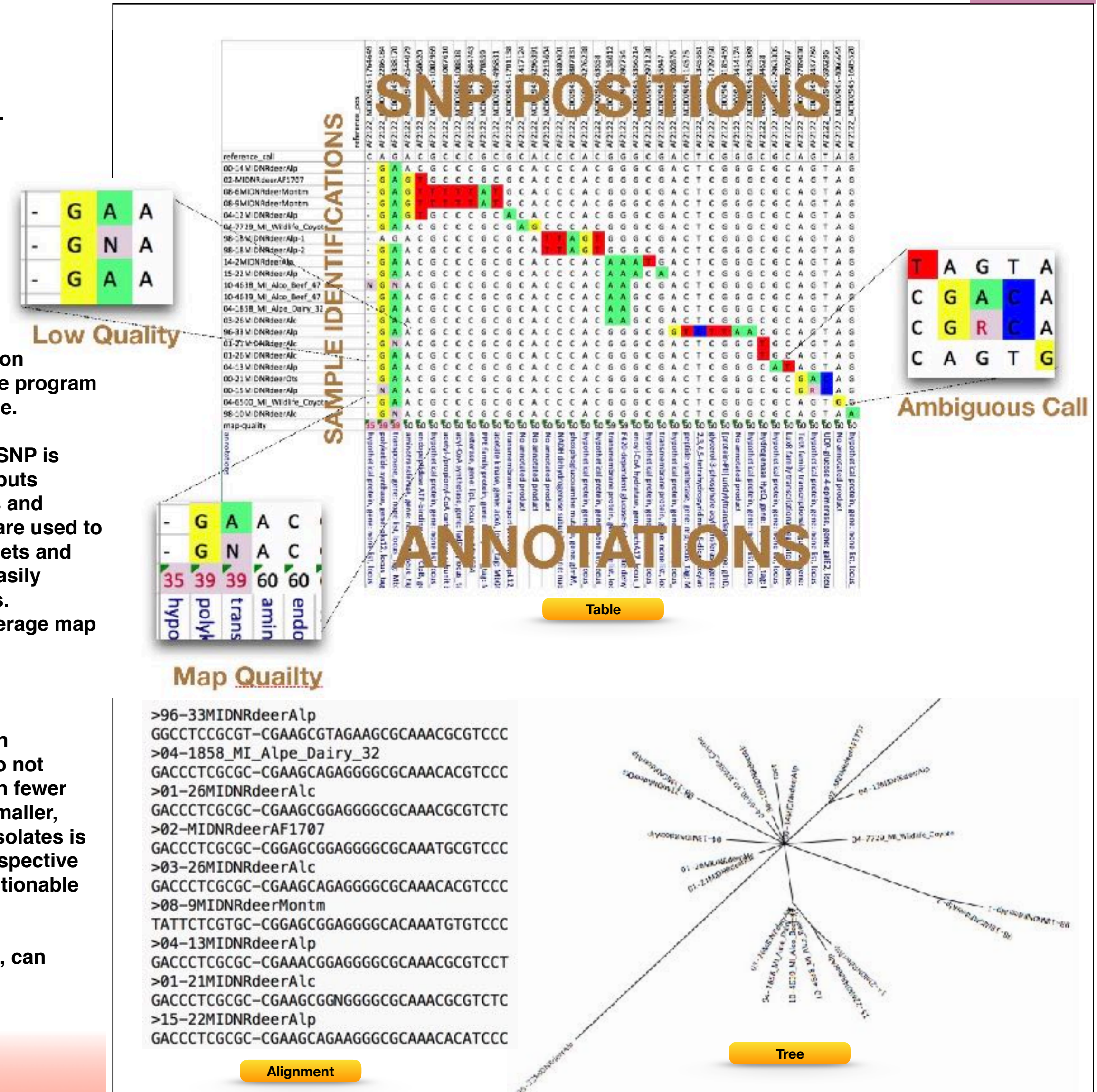
**Pros:**
• Proven pipeline workflow
• Maintain sense of data ownership
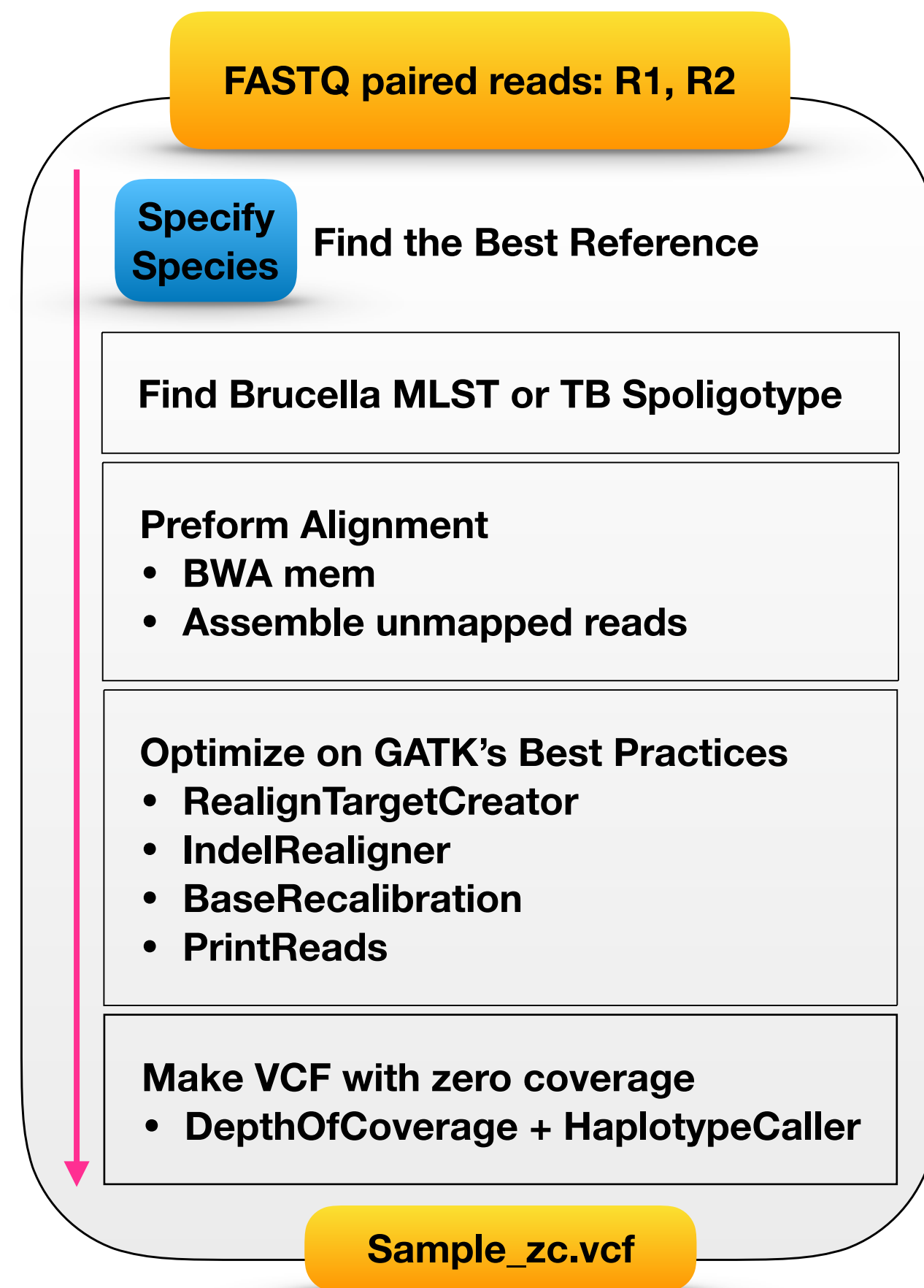• High resolution, validated genotyping

**Cons:**
• Command-line drive
• Unable to run open source tools (GATK) in Windows

## Output



SNP POSITIONS

SAMPLE IDENTIFICATIONS

ANNOTATIONS

Low Quality

Ambiguous Call

Map Quality

Table

```
>96-33MIDNRdeerAlp
GGCCTCCGCGT-CGAAGCGTAGAAGCGCAAACGCGTCCC
>04-1858_MI_Alpe_Dairy_32
GACCCTCGCGC-CGAAGCAGAGGGGCGCAAACACGTCCC
>01-26MIDNRdeerAlc
GACCCTCGCGC-CGAAGCGGAGGGGCGCAAACGCGTCTC
>02-MIDNRdeerAF1707
GACCCTCGCGC-CGGAGCGGAGGGGCGCAAATGCGTCCC
>03-26MIDNRdeerAlc
GACCCTCGCGC-CGAAGCAGAGGGGCGCAAACACGTCCC
>08-9MIDNRdeerMontm
TATTCTCGTGC-CGGAGCGGAGGGGCACAAATGTGTCCC
>04-13MIDNRdeerAlp
GACCCTCGCGC-CGAAACGGAGGGGCGCAAACGCGTCCT
>01-21MIDNRdeerAlc
GACCCTCGCGC-CGAAGCGGNGGGGCGCAAACGCGTCTC
>15-22MIDNRdeerAlp
GACCCTCGCGC-CGAAGCAGAAGGGCGCAAACACATCCC
```

Alignment

Tree

**Mbovis-01D2**

**FASTQ paired reads: R1, R2**

**Specify Species**  **Find the Best Reference**

Find Brucella MLST or TB Spoligotype

Preform Alignment
• BWA mem
• Assemble unmapped reads

Optimize on GATK's Best Practices
• RealignTargetCreator
• IndelRealigner
• BaseRecalibration
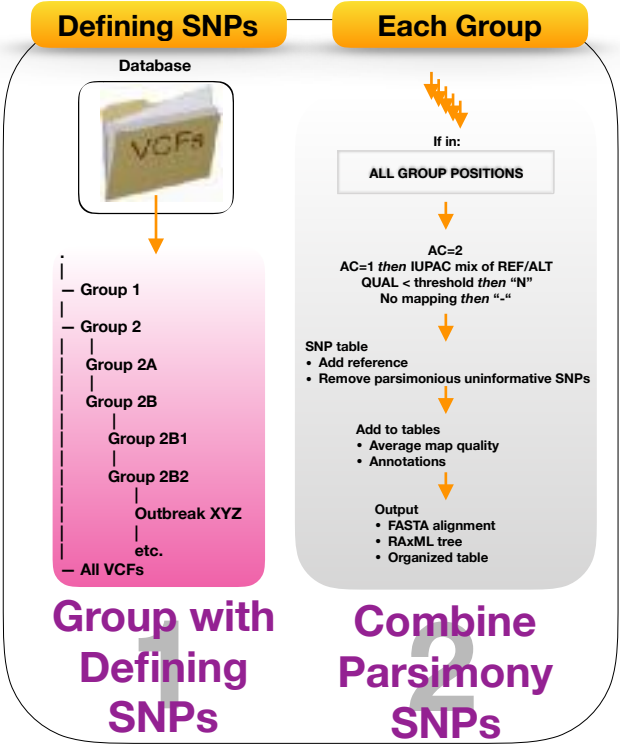• PrintReads

Make VCF with zero coverage
• DepthOfCoverage + HaplotypeCaller

**Sample_zc.vcf**

**FASTQ paired reads: R1, R2**

**Table**

**Defining SNPs**   **Each Group**

Database

VCFs

If in:

ALL GROUP POSITIONS

↓

AC=2
AC=1 then IUPAC mix of REF/ALT
QUAL < threshold then "N"
No mapping then "-"

.
|
— Group 1
|   |
— Group 2
|   |
|   — Group 2A
|   |
|   — Group 2B
|   |   |
|   |   — Group 2B1
|   |   |
|   |   — Group 2B2
|   |   |
|   |   — Outbreak XYZ
|   |   |
|   |   — etc.
|   |
— All VCFs

SNP table
• Add reference
• Remove parsimonious uninformative SNPs

Add to tables
• Average map quality
• Annotations

Output
• FASTA alignment
• RAxML tree
• Organized table

**Group with Defining SNPs**   **Combine Parsimony SNPs**

AC=2, SNP only, QUAL > threshold

**Apply Filters**

ALL QUALITY GROUP POSITIONS

**Find Quality SNPs in each VCF**

Collect Stat Summary
• species used
• number of VCF in comparison
• options list
• run time
• corrupt file list
• ambitious defining SNPs
• groups for each sample
• list of files not renamed

Corrupt VCF check

check duplicates

*Genotyping codes
*Samples to remove

Specie specific parameters
• Qual threshold
• N threshold
• Dependent directory
• Genotyping codes
• Defining SNPs
• Filter regions
• Remove from analysis

-s species

-s species
-e elites
-a all vcfs
-f filter finder

**Figure 2:** Script 2 workflow — blue diamonds represent script options, yellow boxes represent input from dependency files. Those starred are optional.

# Mbovis-01D2