

snakemake project for variant calling

python group

September 5, 2019

variant calling workflow using snakemake language

variant calling is the identification of nucleotide difference on a given reference genome or a transcriptome. Variant calling has become widely accepted in human genetics as a way of identifying variants associated with a specific trait, population or hereditary diseases. Using the standard pipeline available for identification of variant calling from H3ABionet community. We planned to make a more portable and reproducible workflow, for ease of analysis on any given platform. snakemake language offers this opportunity, due to its ability to work across different platform and its use of the python syntax, which is a pipeline language.

The pipeline was designed based on the standard operating procedures (SOPs) from H3Africa website. The pipeline was divided into three phases : phase one: preprocessing of reads (fastq analysis (fastqc), adapter removal and contaminant removal (trimmomatic)), Phase two: Initial variant discovery (alignment to reference genome (bwa and samtools), deduplication (samtools), base quality score recalibration (BSQR protocol from GATK)), Phase three: variant annotation and prioritization (SNP and INDEL variant prioritization (VSQR protocol from GATK)).

METHODOLOGY

To create the snakemake workflow:

Install either anaconda/bioconda platform

First set the environment for analysis: in our repo we already have given the instructions to follow: `setting_up_the_environment`

After setting the environment one can access the snakemake code from : from the repo

```
# install packages (running conda command on r studio)
#install.packages("reticulate")
#install.packages("tidyverse")
# library installation
library(reticulate)
library(tidyverse)

## -- Attaching packages ----- tidyverse 1.2.1 --

## v ggplot2 3.1.0      v purrr   0.2.5
## v tibble  2.1.3      v dplyr  0.8.3
## v tidyr   0.8.2      v stringr 1.3.1
## v readr   1.1.1      v forcats 0.4.0

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()

# include python on the r script command
knitr::knit_engines$set(python = reticulate::eng_python)
```

```

# set the environmnt to the directory with the snakmake file
## Set working directory.
setwd("/home/icipe/Variant_Calling_Project-/pipeline/")

#it is include the environments
conda_list()[[1]][1] %>%
  use_condaenv(required = TRUE)

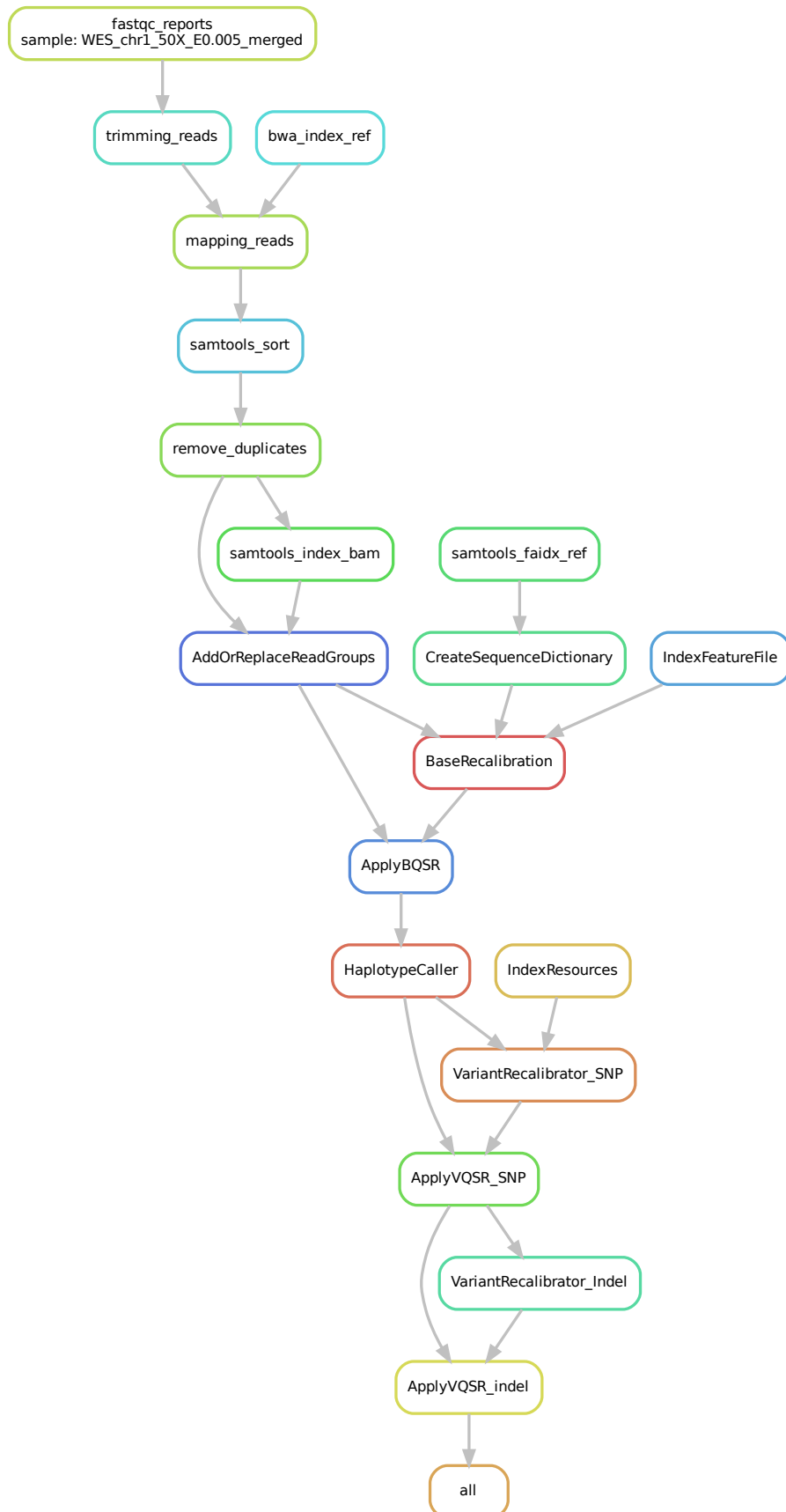
# dry run of the snakemake command
# use of intern = true is used to display knitr output in either pdf or html
system("/home/icipe/miniconda3/envs/variant_calling/bin/snakemake -np")

#Running the command after confrimring with the dry run command

system("/home/icipe/miniconda3/envs/variant_calling/bin/snakemake")

# library installation
library(reticulate)
library(tidyverse)
# displays the workflow directly
system("/home/icipe/miniconda3/envs/variant_calling/bin/snakemake --dag |dot |display " )

```



Snakemake is a diverse language, that can be used for manipulation of data, in this example: we would like to display only phase one of the variant calling analysis.

Employing rmarkdown and commands from snakemake we demonstrate the versatility of the workflow language

```
#Also one has the option to run any number of rules they require
# library installation
library(reticulate)
library(tidyverse)
# command from snakemake (displays the output of phase one script: preprocessing of reads
#(fastqc analysis(fastqc), adapter removal and contaminate removal(trimmomatics)))

# to run a dry run of the snakemake command

# removing intern = true displays the results in the console
system("/home/icipc/miniconda3/envs/variant_calling/bin/snakemake -n --until trimming_reads")

# use of intern = true is used to display knitr output in either pdf or html
system("/home/icipc/miniconda3/envs/variant_calling/bin/snakemake -n --until trimming_reads", intern = '

## [1] ""
## [2] "rule fastqc_reports:"
## [3] "    input: Data/reads/WES_chr1_50X_E0.005_merged_read1.fq.gz, Data/reads/WES_chr1_50X_E0.005_m
## [4] "    output: analyses/fastqc/WES_chr1_50X_E0.005_merged_read1_fastqc.html"
## [5] "    log: logs/fastqc/WES_chr1_50X_E0.005_merged.log"
## [6] "    jobid: 1"
## [7] "    benchmark: benchmarks/fastqc/WES_chr1_50X_E0.005_merged.txt"
## [8] "    wildcards: sample=WES_chr1_50X_E0.005_merged"
## [9] ""
## [10] ""
## [11] "rule trimming_reads:"
## [12] "    input: Data/reads/WES_chr1_50X_E0.005_merged_read2.fq.gz, Data/reads/WES_chr1_50X_E0.005_m
## [13] "    output: analyses/trimmed/WES_chr1_50X_E0.005_merged_read1.paired.fastq.gz, analyses/trimme
## [14] "    log: logs/trimming/WES_chr1_50X_E0.005_merged.log"
## [15] "    jobid: 0"
## [16] "    benchmark: benchmarks/trimming/WES_chr1_50X_E0.005_merged.txt"
## [17] "    wildcards: sample=WES_chr1_50X_E0.005_merged"
## [18] ""
## [19] "Job counts:"
## [20] "\tcount\tjobs"
## [21] "\t1\tfastqc_reports"
## [22] "\t1\ttrimming_reads"
## [23] "\t2"
```

```
#Also one has the option to run any number of rules they require

###library installation
library(reticulate)
library(tidyverse)

# to run a dry run of the snakemake command

# this useful for the application of snakemake command
system("/home/icipc/miniconda3/envs/variant_calling/bin/snakemake -n --help")
```

The snakemake pipeline was evaluated with data used for accreditation in H3Africa consortium. The data was

evaluated to have reported 27921 SNPs and 1589 INDELs demonstrating the functionality of our work flow. Although the snakemake workflow was able to generate the SNPs vcf files. Additional study is required for other types of variant studies.