# Finding *de novo* tandem repeats using VNTRseek

Gisselle Hidalgo[1,2], Gary Benson, Ph.D.[3,4]

[1]Kean University, Union, NJ, [2]Boston University Bioinformatics BRITE REU Program, Summer 2022, Boston, MA
[3]Department of Biology, Boston University, Boston, MA, [4]Department of Computer Science, Boston University, Boston, MA

## Abstract

VNTRseek is a computational tool developed by the Benson Lab to identify sequencing reads containing tandem repeats (TRs) and map them to a set of reference TRs. This project aimed to establish a new computational pipeline to identify potential *de novo* TRs. It comprises Python programs and software tools for processing sequencing data, including SAMtools, BEDtools, and SQLite. With the use of VNTRseek output files as input for my pipeline, these candidate *de novo* TR reads are extracted and aligned to the reference genome and reference set to locate any TR sequences. BEDtools was utilized to determine which reads did not overlap with the reference TRs. The start and end of these reads with our target size TRs were used to locate those sequences in the reference genome. These location ranges were then processed through Tandem Repeat Finder, a TR finder program created by the Benson Lab, to determine whether a TR was found at those locations. Identifying *de novo* TRs can help determine which genomic regions have the propensity for variability and provide a more complete compendium of variable TRs in the human genome.

## Introduction and Background

- Tandem repeats are sequences of two or more DNA bases repeated identically or approximately on a chromosome.
  - Some TR unstable expansions are known to cause genetic diseases such as Huntington's Disease and fragile x syndrome.
- *De novo* TRs are, in this case, where the reference sequence does not indicate a TR.
  - My pipeline will extract these *de novo* TRs and I will perform an alignment to location to determine if there are known TRs at candidate locations.
- VNTRseek is a computational tool designed by the Benson Lab to find tandem repeats (TRs) and map them to a reference set of TRs.
- Variable number tandem repeats (VNTRs) are TR loci that vary in copy number among individuals (shown in Figure 1).



Figure 1 (above): Output from VNTRseek showing VNTRs in an individual.

## Methods

On the right is a figure showing a simplified version of the pipeline created for this project.

- Tools used:
  - BEDtools, SAMtools, VNTRseek, SQLite
- Files types used:
  - SAM
    - Text file format containing alignment information of sequences mapped against reference sequences
  - BAM
    - Binary versions of SAM files
  - BED
    - Tab-delimited files describing features by chromosome, start, end, name, score and strand.
    - DB
      - database files containing fasta reads.
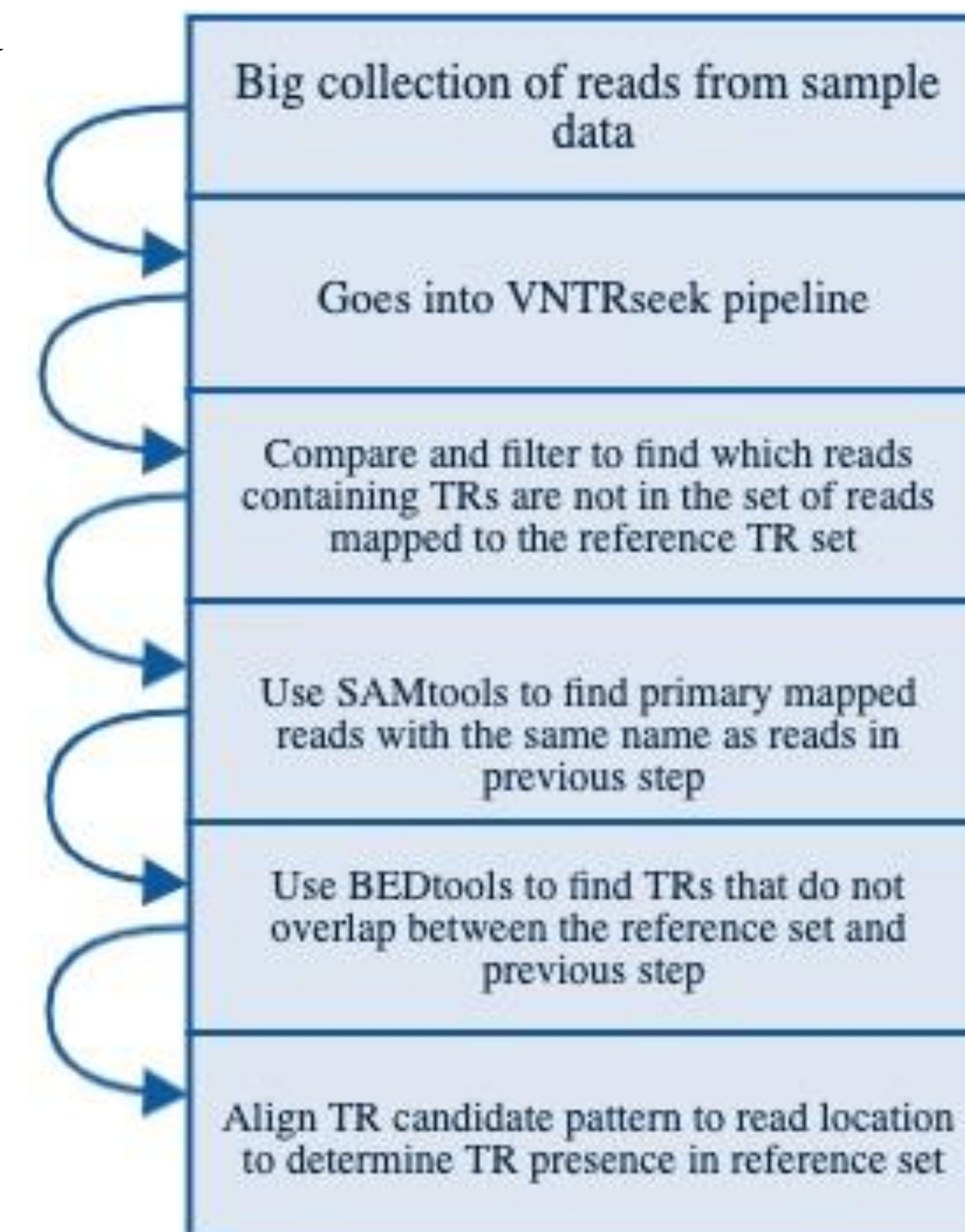  - Files produced from previous step used as input for the next step



Figure 2 (above): Simplified pipeline steps

## Results



Figure 3 (above): Chromosome 1 candidate *de novo* reads.

Preliminary results show over 3,100 candidate *de novo* tandem repeats in bed file format. We focused on chromosome 1 to narrow the results down.

Figure 3 shows the chromosome 1 candidate reads with start and end positions for the repeats.



Figure 4 (above): Highlighted candidate read of interest with corresponding information.

We focused on specific reads with pattern sizes equal to or greater than 20 bp, shown above in Figure 4. Those start and end positions were extracted from the reference genome and run through Tandem Repeats Finder. From the read of interest, it was seen that there was a consensus pattern of CCGTTCC which partially matches the pattern highlighted in Figure 4. This read shows a TR in the reference genome that was not found in the Tandem Repeats Database.



Figure 5 (above): Read of interest with consensus pattern and number of repeats found.

## Conclusions

- The pipeline was successful in filtering which reads were not present in the reference set of TRs.
  - Allowed me to determine which reads could possibly be *de novo*.
  - Had preliminary results of 3,100 reads.
- Found at least one read from the candidate reads bed file and performed an alignment to location to determine if a TR was present at that location.
- It was seen there was a partial match to the original pattern size.
- Many reads chosen for analysis contained TRs not found within the read start and end positions in the reference.
- It is plausible that there are more reads with unmapped TRs in this file in the other chromosomes.
- Further alignment to location will be done to analyze the rest of the candidate reads.

## Future Work

- I would like to look at chromosomes where genetic diseases are commonly found (e.g. X chromosome).
- I would also like to see if this program works with other species, more specifically, *Drosophila*.

## References

M. Rasekh, Y. Hernandez, S. Drinan, J. Fuxman-Bass, and G. Benson (2021). Genome-wide characterization of human minisatellite VNTRs: population-specific alleles and gene expression differences. *Nucleic Acids Research*, 2021, Vol. 49, No. 8, 4308-4324.

Yevgeniy Gelfand, Yozen Hernandez, Joshua Loving, Gary Benson (2014). VNTRseek—a computational tool to detect tandem repeat variants in high-throughput sequencing data. *Nucleic Acids Research*, Vol. 42, No. 14, 8884–8894.

G. Benson,"Tandem repeats finder: a program to analyze DNA sequences" *Nucleic Acids Research* (1999) Vol. 27, No. 2, pp. 573-580.

Yevgeniy Gelfand, Alfredo Rodriguez and Gary Benson TRDB—The Tandem Repeats Database. (2007) *Nucleic Acids Research*, 35:D80-D87.

## Acknowledgement