# *Arabidopsis thaliana,* Thale Cress, Genome Assembly

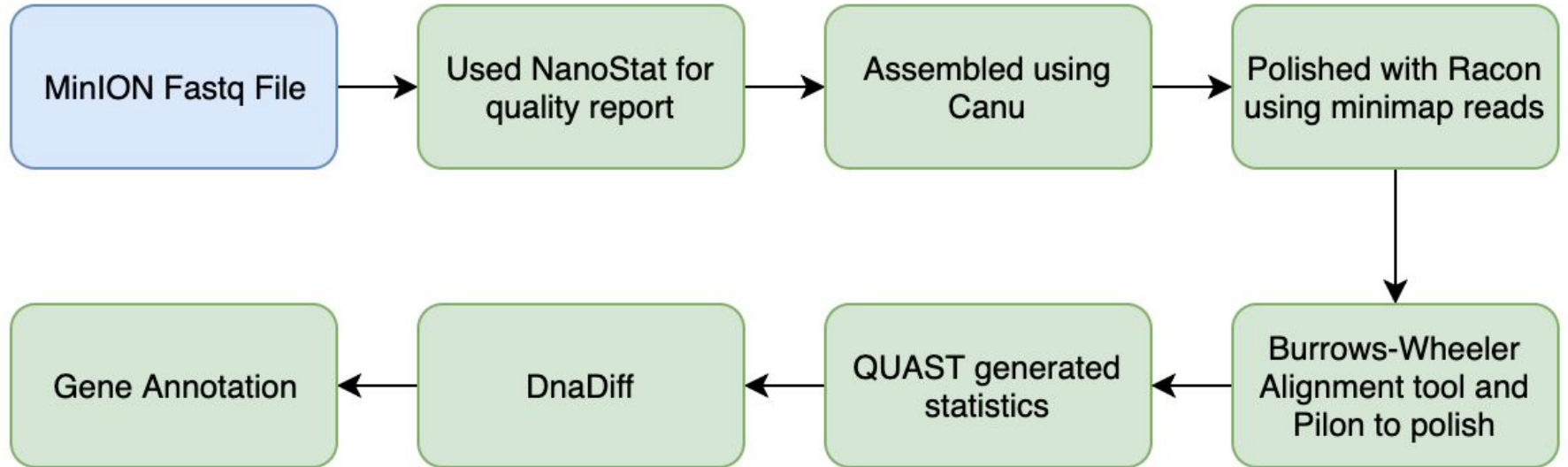By Jordan Callahan

# Purpose of Project

-   Assemble complete genome using only
    minION data
-   Establish high contiguity genome with
    racon  and pilon polishing
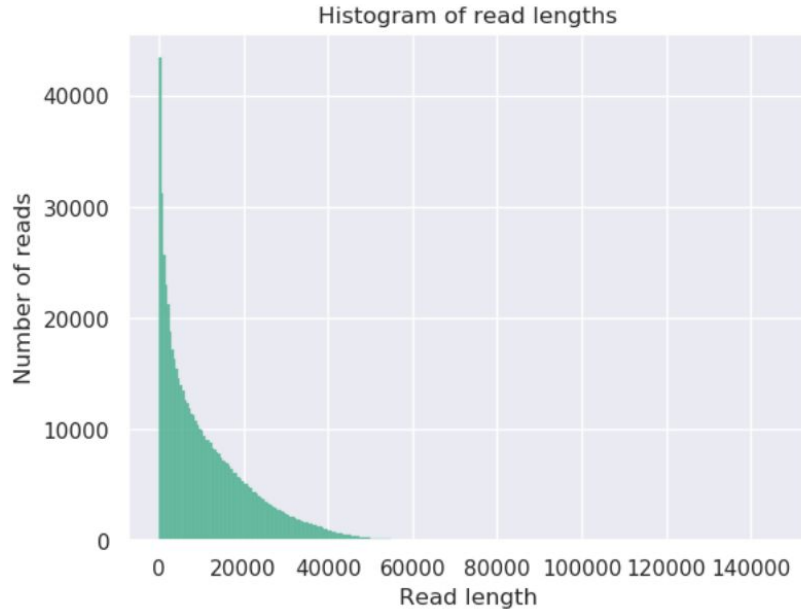-   Observe the top hits for protein-encoding
    genes

# Workflow

# Nano Plot

- 309947 reads of the original ~9 million did not pass quality score of 7

Histogram of read lengths



## Summary statistics

**General summary:**

**Mean read length: 11,417.1**

**Mean read quality: 0.0**
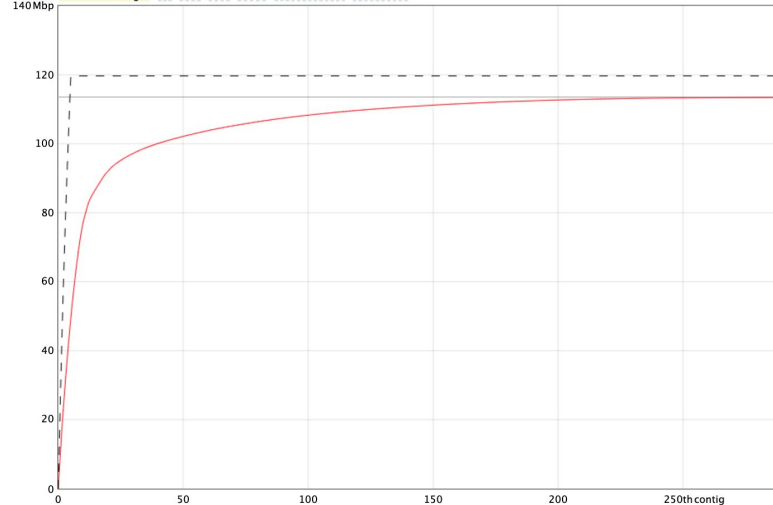
**Median read length: 8,208.0**

**Median read quality: 0.0**

**Number of reads: 613,080.0**

**Read length N50: 20,006.0**
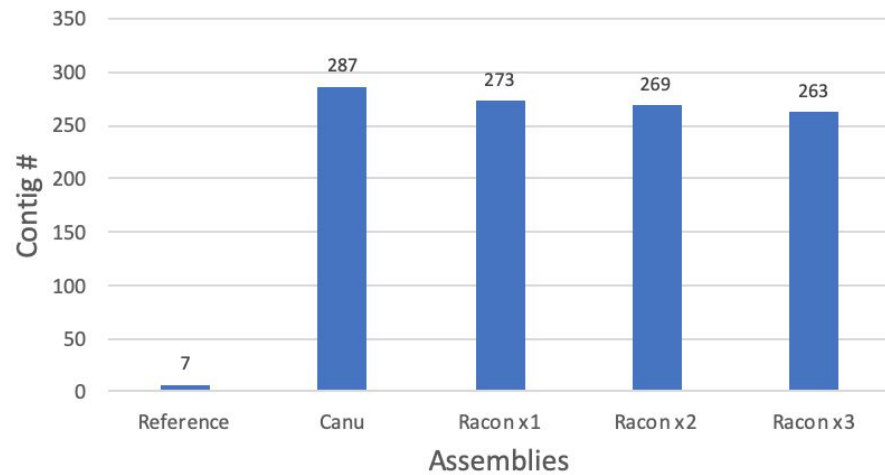
**Total bases: 6,999,607,984.0**

# QUAST

# Assembly Size

## Assembly size with each round of polishing
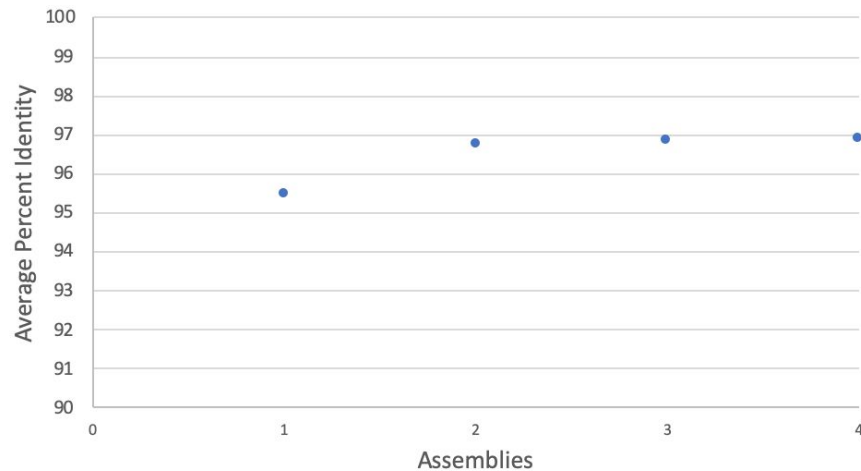


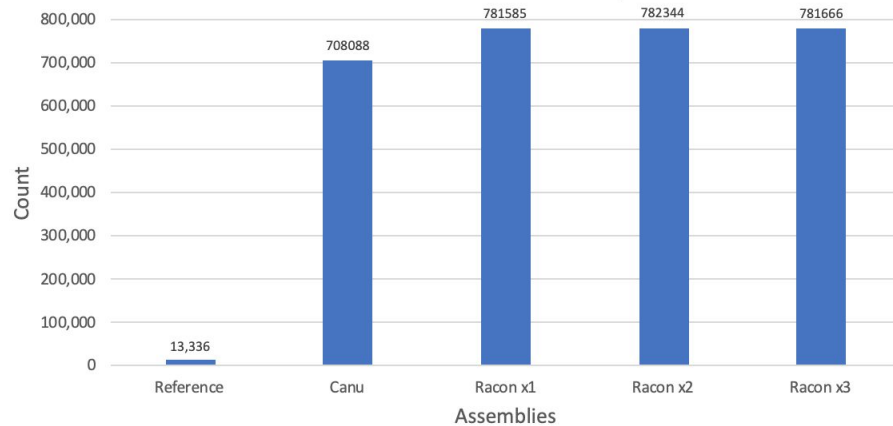- Reference assembly size is 119,667,750

# DnaDiff report



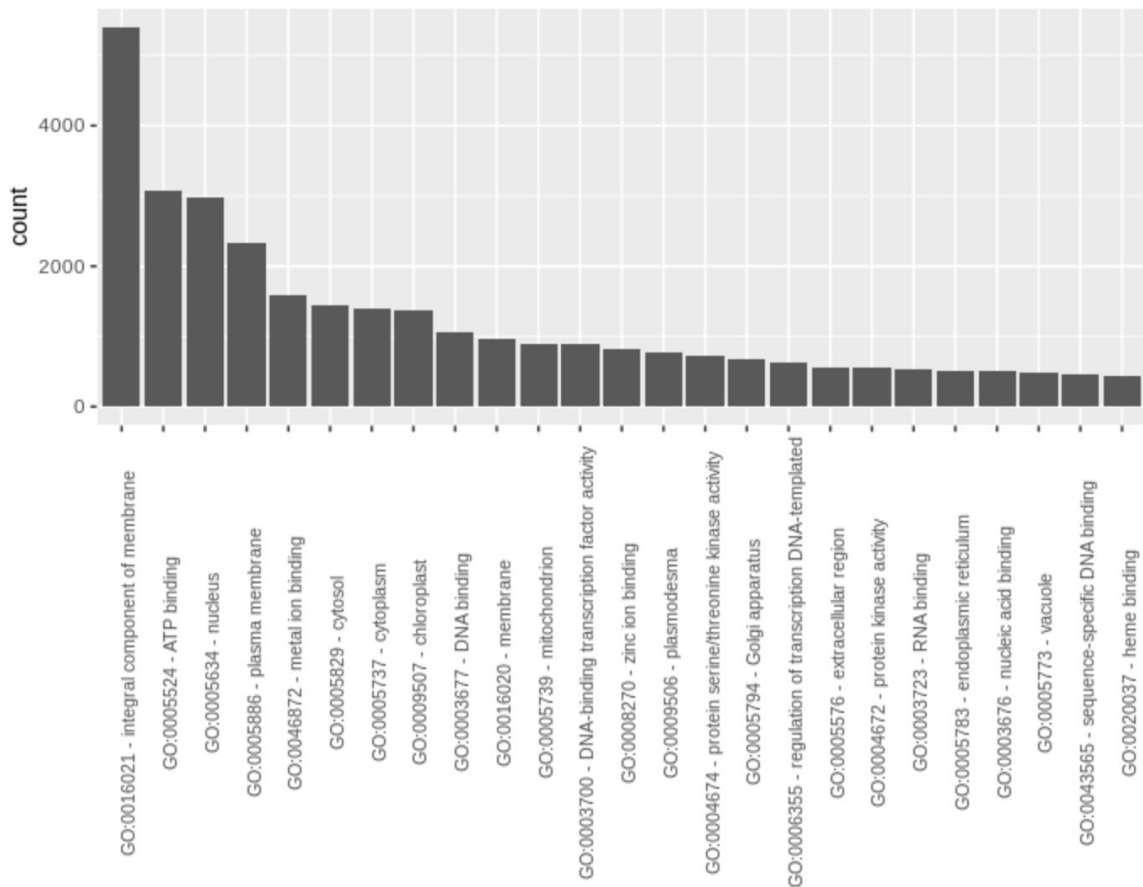Average Percent ID to Reference Genome

1 - Canu   2 - Racon x1   3 - Racon x2   4 - Racon x3

Number of SNPs in each assembly

Top 25 gene hits

# Gene Annotation

- Prodigal gene prediction to create a list of proteins
- GOFEAT website
- Read csv file into ggplot
- Top 25 gene hit

# Concluding Statements

- A complete genome was successfully assembled and analyzed

- Workflow was altered from original experiment (Miniasm, blastn)

- From the paper, researchers were able to assemble into 62 contigs with an N50 of 12.3 Mb.

# References

Michael, T.P., Jupe, F., Bemm, F., Motley, S.T., Sandoval, J.P., Lanz, C., Loudet, O., Weigel, D.

and Ecker, J.R., 2018. High contiguity Arabidopsis thaliana genome assembly with a single nanopore flow cell.

*Nature communications*, *9*(1), p.541.