# LABORATORY REPORT #5
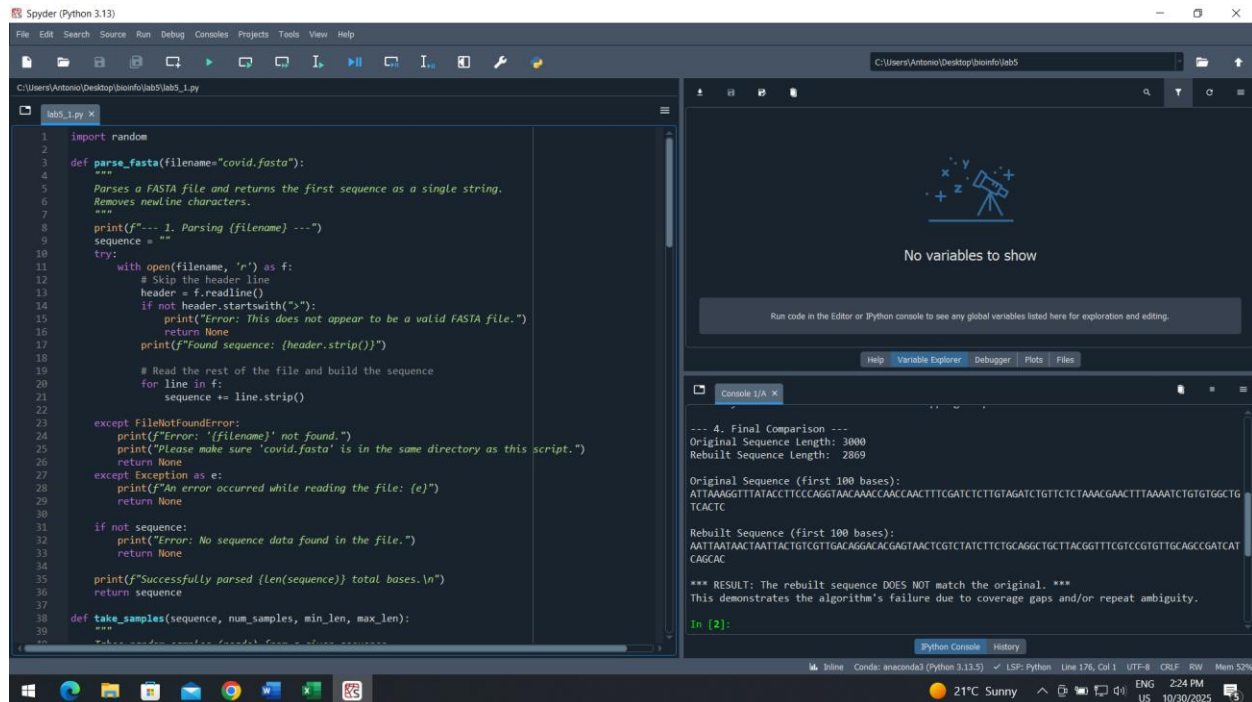
GEORGESCU Mihai-Antonio, 1242EA

Bioinformatics, 4th year 1st semester, 2025-2026

lab5_1.py



--- 1. Parsing covid.fasta ---

Found sequence: >NC_045512.2 Severe acute respiratory syndrome coronavirus 2 isolate Wuhan-Hu-1, complete genome

Successfully parsed 29903 total bases.

--- 2. Taking 2000 random samples ---

Generated 2000 samples (e.g., 'GAAGTGTTAACAGAGGAAGTTGTCTTGAAAACTGGTGATTTACAACCATTAGAACAACCTACTAGTGAAGCTGTTGAAGCTCCATTGGTTGGTACACCAGT...').

--- 3. Attempting to rebuild sequence (min_overlap=10) ---

Assembly stuck. No more forward overlaps found.

Assembly finished. Found 411 total overlapping samples.

--- 4. Final Comparison ---

Original Sequence Length: 3000

Rebuilt Sequence Length:  2869

Original Sequence (first 100 bases):

ATTAAAGGTTTATACCTTCCCAGGTAACAAACCAACCAACTTTCGATCTCTTGTAGATC
TGTTCTCTAAACGAACTTTAAAATCTGTGTGGCTGTCACTC

Rebuilt Sequence (first 100 bases):

AATTAATAACTAATTACTGTCGTTGACAGGACACGAGTAACTCGTCTATCTTCTGCAG
GCTGCTTACGGTTTCGTCCGTGTTGCAGCCGATCATCAGCAC

*** RESULT: The rebuilt sequence DOES NOT match the original. ***

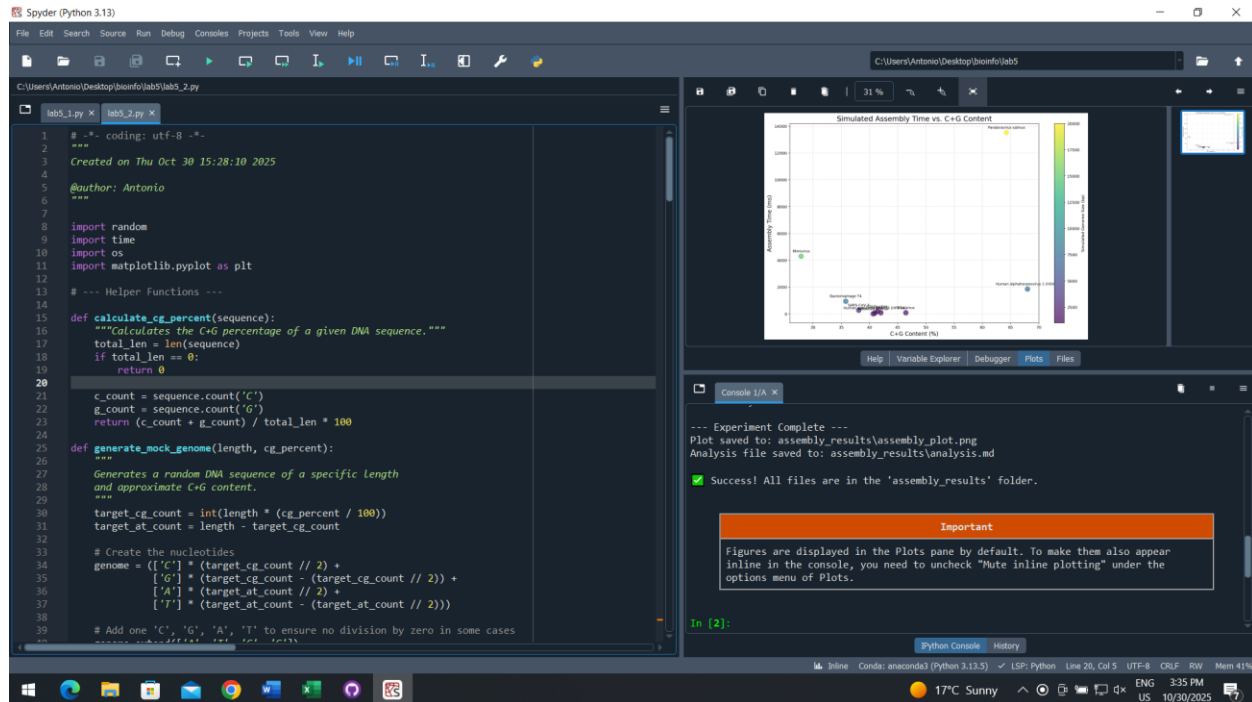This demonstrates the algorithm's failure due to coverage gaps and/or repeat ambiguity.

The main problem with this algorithm is that it's local and it makes decisions based only on the immediate, best-looking overlap it can find, without any knowledge of the larger, global structure. This approach fails because it has no way to resolve ambiguity.

The 10-base minimum overlap is just a hint, but in a large sequence, this 10-base string can appear coincidentally in many places, or it can appear legitimately in many places (a repeat). The algorithm can't tell the difference and will either:

1. Get stuck and halt the assembly, as my script showed (a coverage gap or an unresolvable fork).

2. Mis-assemble the sequence by connecting two fragments that shouldn't be together.

The specific structures in the original sequence that cause these computational issues are all forms of repeats, so the algorithm fails because the samples are shorter than the repeats. The 10-base overlap is just a rule for a potential connection, but it's not nearly enough information to prove that a connection is correct.

lab5_2.py



--- Starting Assembly Experiment ---

Processing: Human papillomavirus 16 (HPV-16) (Simulated length: 1000 bp)
   Generated 64 samples.
   Assembly Time: 64.42 ms

Processing: HIV-1 (Simulated length: 1200 bp)
   Generated 76 samples.
   Assembly Time: 80.40 ms

Processing: Zika virus (Simulated length: 1300 bp)
   Generated 83 samples.
   Assembly Time: 81.25 ms

Processing: Influenza A virus (Simulated length: 1500 bp)
   Generated 96 samples.
   Assembly Time: 3.20 ms

Processing: Ebola virus (Simulated length: 2000 bp)
  Generated 128 samples.
  Assembly Time: 179.95 ms

Processing: SARS-CoV-2 (Simulated length: 3000 bp)
  Generated 192 samples.
  Assembly Time: 270.07 ms

Processing: Human alphaherpesvirus 1 (HSV-1) (Simulated length: 7500 bp)
  Generated 480 samples.
  Assembly Time: 1848.38 ms

Processing: Bacteriophage T4 (Simulated length: 8000 bp)
  Generated 512 samples.
  Assembly Time: 934.13 ms

Processing: Mimivirus (Simulated length: 15000 bp)
  Generated 960 samples.
  Assembly Time: 4291.91 ms

Processing: Pandoravirus salinus (Simulated length: 20000 bp)
  Generated 1280 samples.
  Assembly Time: 13516.29 ms

--- Experiment Complete ---

Simulated Assembly Time vs. C+G Content