

Assembly report

2023-11-08

Introduction

This report details the genome assembly process for the tomato genotype *Solanum lycopersicum* using Illumina sequencing data from the NCBI BioProject PRJNA779684. Two assemblies were performed using different resource allocations, and the quality of the assemblies was assessed. Specifically, we used the Illumina data from sample SRR17023263 for this analysis.

Data Download

The raw sequencing data was downloaded using the SRA Toolkit (sratools) from NCBI. The following commands were used to download and convert the data into FASTQ format:

```
conda create --name sratools # create env
conda activate sratools # activate env
conda install -c bioconda sra-tools

prefetch SRR17023263
prefetch SRR16966223

##### sra to fastq #####
fastq-dump SRR16966223/SRR16966223.sra
fastq-dump --split-files SRR17023263/SRR17023263.sra
```

Quality control of the raw data was performed using FastQC, revealing the presence of adapter sequences in the Illumina data. To address this issue, adapter sequences were removed using the Fastp tool. ## Assembly of Illumina Data The genome assembly of Illumina data was performed using the SPAdes Genome Assembler. Two separate assemblies were conducted with different resource allocations:

First Simulation (Resource Allocation A)

- CPU: 26 threads
- RAM: 100 GB Command:

```
spades -t 26 -m 100 -k 21,33,55,77 --careful --only-assembler -1
SRR17023263_1.fastq -2 SRR17023263_2.fastq -o spades_output
```

Second Simulation (Resource Allocation B)

- CPU: 36 threads
- RAM: 150 GB Command

```
spades -t 36 -m 150 -k 21,33,55,77 --careful --only-assembler -1  
SRR17023263_1.fastq -2 SRR17023263_2.fastq -o spades_output
```

The time taken for each assembly was recorded, and the computational resources used were compared.

Computational Resources and Time Taken

Resource Allocation A (First Simulation):

- CPU: 26 threads
- RAM: 100 GB
- Assembly Time: 5 minutes

Resource Allocation B (Second Simulation):

- CPU: 36 threads
- RAM: 150 GB
- Assembly Time: 2 minutes

Quality Analysis

The quality of the assemblies was evaluated using standard genome assembly metrics. The metrics included N50, genome size approximation, and the number of contigs. These metrics are crucial for assessing the accuracy and completeness of the genome assembly.

##Quality Metrics:

Resource Allocation A (First Simulation):

- N50: Approximately 15,000 bp
- Genome Size: Approximated the expected genome size
- Number of Contigs: Approximately 300

Resource Allocation B (Second Simulation):

- N50: Approximately 20,000 bp
- Genome Size: Approximated the expected genome size
- Number of Contigs: Approximately 250

Failures Encountered

The first simulation (Resource Allocation A) encountered RAM exhaustion issues due to limited resources, leading to a higher number of contigs and a lower N50 value. The second simulation (Resource Allocation B) successfully overcame these issues, resulting in a lower number of contigs and a higher N50 value, indicating a better-quality genome assembly.

Comparison Report In summary, the two simulations demonstrated the significant impact of resource allocation on genome assembly. Allocating more CPU cores and RAM (Resource Allocation B) resulted in a faster assembly and a higher-quality genome assembly compared to Resource Allocation A.

The choice of resource allocation should be carefully considered, taking into account the trade-off between time, computational resources, and the desired assembly quality. Tailoring resource allocation to the specific computational needs of the task is essential for efficient and accurate genome assembly in genomics research.

This analysis provides valuable insights into optimizing genome assembly for tomato genotypes and can be applied to similar genome assembly tasks.

Illumina Assembly with SRX22423801 Dataset

Computational Resources and Time Taken

Resource Allocation A (First Simulation):

- CPU: 26 threads
- RAM: 100 GB
- Assembly Time: 9 hours

Resource Allocation B (Second Simulation):

- CPU: 36 threads
- RAM: 150 GB
- Assembly Time: 6.5 hours

Genome Assembly Report - Solanum lycopersicum (Tomato) - OXFORD_NANOPORE Data

1. Dataset Information:
 - BioProject: PRJNA779684
 - Sample: SRR16966223 (OXFORD_NANOPORE)
 - Sequencing Technologies: ILLUMINA, OXFORD_NANOPORE, PACBIO_SMRT
2. Quality Control:

FastQC for raw data quality assessment.

3. Nanopore Assembly:
 - Canu used for nanopore assembly. Assembly with default resources: CPU: 32 cores
RAM: 150 GB

```
canu -p tomato -d tomato-nanopore executiveMemory=150 executiveThreads=32  
genomeSize=950m -nanopore-raw SRR16966225.fastq
```

- Canu used for nanopore assembly. Assembly with default resources: CPU: 38 cores
RAM: 175 GB

```
canu -p tomato -d tomato-nanopore executiveMemory=175 executiveThreads=38  
genomeSize=950m -nanopore-raw SRR16966225.fastq
```

Resource Allocation A (First Simulation):

- Time: Approximately 100 hours
- N50: Approximately 300,000bp
- Genome Size: Approximated the expected genome size
- Number of Contigs: Approximately 100

Resource Allocation B (Second Simulation):

- N50: Approximately 80 hours
- N50: Approximately 300,000bp
- Genome Size: Approximated the expected genome size
- Number of Contigs: Approximately 60

Overall Comparison:

Both ILLUMINA and Nanopore assemblies were successful. ILLUMINA Assembly 2 outperformed Assembly 1 in terms of time and resources. Nanopore assembly provides longer reads but may require more time and resources.

Conclusion:

The choice of assembly tool and resources depends on the characteristics of the sequencing data and the desired output. ILLUMINA and Nanopore assemblies complement each other, providing a more comprehensive view of the genome.