

DESIGN AND ANALYSIS OF ALGORITHMS PROJECT

DE Bruijn Algorithm Implementation

Group Members

Ali Amar	426014
Ibrahim Qaiser	405459

Table of Contents:

1. Project Description
2. Background
3. Research Goals
4. Risks/Ethical Concerns

de Bruijn Graph Algorithm for Genome Sequencing

Project Description

Problem Statement:

The problem of genome sequencing is a fundamental challenge in bioinformatics. It involves determining the precise order of nucleotides in a DNA sequence. Traditional sequencing methods are time-consuming and expensive, making them impractical for large-scale sequencing projects. Therefore, there is a need for efficient algorithms that can accurately reconstruct the genome from short DNA fragments.

Different Solutions to the Problem:

Several approaches have been proposed to address the challenge of genome sequencing, including overlap-layout-consensus, shotgun sequencing, and de novo assembly. Among these, the de Bruijn graph algorithm stands out as a powerful and widely used method.

The de Bruijn Graph Algorithm:

The de Bruijn graph algorithm is a graph-based approach that constructs a de Bruijn graph from the given short DNA fragments. This graph represents the overlaps between the fragments by breaking them into k-Mers, where k is a fixed length. Each k-Mer is represented as a node in the graph, and an edge is created between two nodes if there is an overlap of length k-1 between their respective k-Mers.

The de Bruijn graph provides a compact representation of the DNA fragments and their overlaps, enabling efficient assembly of the genome. By traversing the graph, the algorithm can reconstruct the entire genome sequence by finding a Eulerian path, which visits each edge exactly once.

Advantages:

Compared to other sequencing methods, the de Bruijn graph algorithm offers several advantages:

- 1. Memory Efficiency:** The de Bruijn graph representation requires less memory compared to storing the entire set of DNA fragments. This makes it feasible to handle large-scale sequencing projects and reduces the computational resources required.

2. **Scalability:** The algorithm scales well with increasing genome size and sequencing depth. It can efficiently handle complex genomes with repetitive regions, which pose challenges for other sequencing methods.
3. **Error Correction:** The de Bruijn graph can incorporate error correction techniques, allowing for accurate assembly even in the presence of sequencing errors or variations in the DNA sequence.
4. **Parallelization:** The algorithm can be parallelized to take advantage of multi-core processors and distributed computing systems, further enhancing its speed and scalability.

By leveraging the strengths of the de Bruijn graph algorithm, our project aims to develop an efficient and accurate genome sequencing tool. We will focus on implementing the algorithm, optimizing its performance, and evaluating its effectiveness using real-world datasets.

Background

Genome assembly represents the process of putting a large number of short DNA sequences back together to recreate the original chromosomes from which the DNA originated. Genome assembly may be divided into two categories: mapping to a reference genome (also called reference-based alignment) and de novo assembly. The process of assembling a new genome from scratch without the use of reference genomic material is known as "de novo assembly." A digital nucleic acid sequence database that serves as a representative representation of a species' gene set is called a reference genome or reference assembly.

De novo whole-genome assembly is compared to putting together a massive, multimillion-piece jigsaw puzzle. Fragmented sequences need to be assembled; whole genome sequences cannot be read all at once. The construction of the roughly 3 billion base pair human genome necessitates billions of small sequences. De novo assembly is hampered by the complexity and non-randomness of genomic sequences. Nonrandom repeat sequences make up around half of the human genome; these elements can lead to assembly errors or gaps. Repeated sequences can cause nonuniform read depth, which might cause gain or loss of copies during assembly.

There are some factors which can heavily affect genome sequencing:

1. Genome properties

- Genome size influences data demand, calculated based on genome sizes of similar species.
- Repetition affects assembly and may result in mis-assemblies. High heterozygosity causes assemblies to become fragmented. Heterozygosity problems are avoided by sequencing haploid tissue.
- Illumina sequencing encounters difficulties due to heterogeneous GC content.

2. Nucleic Acid Extraction

- Quantity, purity, and integrity of DNA/RNA are crucial.
- Excellent nucleic acid purity is necessary for de novo sequencing.
- Chemical purity and structural integrity are significant quality requirements.

3. Sequencing techniques

- Decisions have an impact on success and cost. The use of next-generation sequencing (NGS) in noteworthy initiatives.
- Third-generation sequencing shows potential in extending coverage across repetitive areas.

4. Raw data processing

- Quality-trimmed or raw data may be preferred by assembly tools.
- It's crucial to filter and trim low-quality readings and read endings.

Benefits

Following are some of the benefits of genome assembly as well which also cater the applications in real world:

Medical Genetics: Identifies genetic variations linked to diseases, aiding personalized medicine and precision oncology treatments.

Evolutionary Biology: Enables comparative analysis to study species relationships, identify genomic changes driving evolution, and understand adaptation.

Agricultural Genomics: Characterizes crop genomes to find genes for desirable traits, facilitating breeding programs for improved crops.

Microbial Genomics: Studies diversity and function of microbial communities, enabling discovery of novel genes and biosynthetic pathways.

Environmental Genomics: Analysis complex microbial communities in various environments, aiding conservation efforts.

Human Population Genomics: Investigates genetic diversity, population structure, and ancestry in human populations.

Functional Genomics: Provides a foundation for functional annotation and multi-omics studies to understand gene function, regulatory networks, and complex biological processes.

Future

Following are some of the future points of genome assembly:

1. **Longer Reads:** Longer and more accurate reads will be produced by new sequencing technologies, which will enhance the assembly of repetitive regions and variants.
2. **Merging Technology:** By utilizing their respective advantages, integrating data from several sequencing technologies can improve assembly quality.
3. **Better Algorithms & Software:** Error correction, scaffolding, and assembly polishing will all be enhanced by the ongoing development of specialized algorithms and software. The role of machine learning might increase.
4. **Reference-Free Assembly:** Complex or innovative genome assembly will require new techniques that do not rely on preexisting genomes.
5. **Population & Metagenomics Assembly:** As sequencing costs decline, strategies to address issues such as genetic variation will be needed to assemble genomes from sizable populations and microbiological communities.
6. **Omics Integration:** Deeper understanding of gene function can be obtained by combining genome data with other omics data.

Research Plan Goals

Following are some of the main points around which our research would revolve:

Algorithm Optimization: To tackle repeated sequences, a significant obstacle in genome assembly, we will incorporate machine learning in future exactly not in

initial version of project. Machine learning can enhance polishing, scaffolding, and mistake correction, resulting in more precise and comprehensive assemblies.

Reference-Free Assembly: Our main focus will be on techniques for reference-free assembly of unique and varied genomes. Reference-guided methods perform well on known organisms but have difficulties when dealing with ancient DNA or non-model species. Reconstruction is possible via reference-free assembly, which does not rely on preexisting references. Our goal in optimizing these algorithms is to increase the breadth of genetic research.

Our first effort will be focused on these two areas, providing a foundation for addressing more general genome assembly and analysis difficulties.

Risk/Ethical Concerns

Since the research is computational in nature, neither live things nor non-living things will be specifically impacted. Only the assembly of shorter DNA sequences into longer, continuous DNA sequences or even whole genomes will be simulated by us. In no way is the DE Bruijn graph technique hazardous or risky.

References

<https://www.cd-genomics.com/an-overview-of-genome-assembly.html>

<https://academic.oup.com/bib/article/19/1/23/2339783>

<https://data-science-sequencing.github.io/Win2018/lectures/lecture7/>