

Genome assembly composition of the String "ACGT" array: A review of data Structure accuracy and performance challenges

Sherif Magdy Mohamed Abdelaziz Barakat^{Corresp., 1}, Roselina Sallehuiddin¹, Siti Sophiayati Yuhani², Raja Farhana R. Khairuddin³, Yasir Mahmood⁴

¹ Computer Science, School of Computing, Faculty of Engineering, Universiti Teknologi Malaysia, Skudai, Johor, Malaysia

² Advanced Informatics Department, Razak Faculty of Technology and Informatics, Universiti Teknologi Malaysia, Kuala Lumpur, Kuala Lumpur, Malaysia

³ Department of Biology, Universiti Pendidikan Sultan Idris Tanjung, Malim, Malim, Malaysia

⁴ Faculty of Information Technology, The University of Lahore, Lahore, Lahore, Pakistan

Corresponding Author: Sherif Magdy Mohamed Abdelaziz Barakat
Email address: sherifmagdy.barakat@gmail.com

Background. In this era, biologists have been increasingly interested in interpreting the genome sequence. Unfortunately, the current sequencing technology is limited to providing the genome as a massive number of short strings (reads). Computer algorithms for genome assembly construct the entire genome from reads in two approaches. De novo approach concatenates the reads based on the exact match between their suffix-prefix (overlapping). Reference-guided approach orders the reads based on their offsets in a well-known reference genome (reads alignment). Repetitive sequences in short reads (repeats) present technical ambiguity, making the algorithm unable to distinguish the reads and assemble them incorrectly (misassembly). On the other hand, the massive number of reads causes a big assembly performance challenge. **Method.** The repeat identification method was introduced for misassembly by prior identification of repetitive sequences, creating a repeat knowledge base to reduce ambiguity during the assembly process, thus enhancing the accuracy of the assembled genome. Also, hybridization between assembly approaches resulted in a lower misassembly degree with the aid of the reference genome if it exists. On the other hand, assembly performance is optimized through data structure indexing or parallelization. This paper's primary aim and contribution are to support the researchers through an extensive review to ease other researchers' search for genome assembly studies. We highlighted the most recent developments and limitations in genome assembly accuracy and performance optimization. **Results.** As a result of this review, we found that most repeat identification methods can detect only a specific length of the repeat, while repeat in biology varies. On the other hand, hybridization between de novo and reference-guided approach lacks any repeat resolution. In addition, it is too time-intensive. Regardless of the tremendous work on genome assembly performance optimization, no current method fits the optimization of the hybrid assembly approach. **Conclusion.** We suggest combining multiple repeat identification methods to enhance the accuracy of identifying the repeats as an initial step to the assembly hybrid approach and combining genome indexing with parallelization for better optimization of its performance.

Genome Assembly Composition of the String "ACGT" Array: A Review of Data Structure Accuracy and Performance Challenges

Sherif Magdy Mohamed Abdelaziz Barakat¹, Roselina Sallehuddin¹, Siti Sophiayati Yuhaniz², Raja Farhana R. Khairuddin³, Yasir Mahmood²

¹ Computing Department, School of Computing- Faculty of Engineering, Universiti Teknologi Malaysia Skudai, Joho, Malaysia

² Advanced Informatics Department, Razak Faculty of Technology and Informatics- Universiti Teknologi Malaysia, Kuala Lumpur, Malaysia

³ Department of Biology, Universiti Pendidikan Sultan Idris Tanjung, Malim, Malaysia

Corresponding Author:

Sherif Barakat¹

UTM, Jalan Iman, 81310 Skudai, Johor

Email address: sherifmagdy.barakat@gmail.com

Abstract

Background. The development of sequencing technology increases the number of genomes being sequenced. However, obtaining a quality genome sequence remains a challenge in genome assembly by assembling a massive number of short strings (reads) with the presence of repetitive sequences (repeats). Computer algorithms for genome assembly construct the entire genome from reads in two approaches. The *de novo* approach concatenates the reads based on the exact match between their suffix-prefix (overlapping). Reference-guided approach orders the reads based on their offsets in a well-known reference genome (reads alignment). The presence of repeats extends the technical ambiguity, making the algorithm unable to distinguish the reads resulting in misassembly and affecting the assembly approach accuracy. On the other hand, the massive number of reads causes a big assembly performance challenge.

Method. The repeat identification method was introduced for misassembly by prior identification of repetitive sequences, creating a repeat knowledge base to reduce ambiguity during the assembly process, thus enhancing the accuracy of the assembled genome. Also, hybridization between assembly approaches resulted in a lower misassembly degree with the aid of the reference genome. The assembly performance is optimized through data structure indexing and parallelization. This paper's primary aim and contribution are to support the researchers through an extensive review to ease other researchers' search for genome assembly studies. The study also highlighted the most recent developments and limitations in genome assembly accuracy and performance optimization.

Results. Our findings show the limitations of the repeat identification methods available, which only allow to detect of specific lengths of the repeat, and may not perform well when various

types of repeats are present in a genome. We also found that most of the hybrid assembly approaches, either starting with *de novo* or reference-guided, have some limitations in handling repetitive sequences as it is more computationally costly and time intensive. Although the hybrid approach was found to outperform individual assembly approaches, optimizing its performance remains a challenge. Also, the usage of parallelization in overlapping and reads alignment for genome assembly is yet to be fully implemented in the hybrid assembly approach.

Conclusion. We suggest combining multiple repeat identification methods to enhance the accuracy of identifying the repeats as an initial step to the hybrid assembly approach and combining genome indexing with parallelization for better optimization of its performance.

Introduction

Repeats in the genome

Deoxyribonucleic acid (DNA) is the genetic material of most organisms, which is made up of a chain of four chemical bases indicated by letters A, C, G, and T, and a complete set of DNA sequences an organism called genome [45]. Repetitive sequences were found across all kingdoms of life. More than 50% of the human genome is occupied by DNA repetitive sequences. Repeat is a segment of sequence that appears multiple times in the genome in the identical or near-identical form [14],[47]. The essential repeat categories in biology are transposable elements (TE), and tandem repeats (TR) [1]. TEs are DNA sequences that are able to copy themselves from one genome region, overwriting another region, which are considered to play a major role in the genome evolution, thus, can change the structure and size of genomes, as shown in Figure.1. The length of TE can be varied, such as long-terminal repeat retrotransposons, the length of which generally ranges from 100 bp to 25 kb [64]. On the other hand, TR is a repetitive DNA sequence, which is either known as a microsatellite for a short repeat segment (1- 12 bp) or a minisatellite for a longer repeat segment (12 - 500 bp) located adjacent to each other. Both types of repeats have the ability to expand their copy number and change the structure and size of the genome [2],[3]. The expansion of repeat sequences can cause many diseases in humans, such as Huntington disease and Kennedy disease [4]. The presence of the repetitive offers crucial biological information that stores historical genome changes within and between species, which should be properly handled during assembly, and cannot be treated as a duplication and then truncated by data cleaning.

Genome Assembly in Genome Analysis Project

Genome analysis can be divided into three main phases; genome sequencing, genome assembly, and genome annotation. In the first phase, genomic sequences from an organism will be extracted and further processed through a sequencing machine, in which the sequences will be fragmented into smaller fragments as an output (reads). In the genome assembly phase, all the reads will be overlapped or mapped to the existing reference genome to construct the genome of

the organism. Finally, the assembled genome will be annotated according to their genome position for coding (e.g. gene) and non-coding (e.g. intron) regions. The genome assembly phase is the most computationally intensive part of genome analysis. Given such a collection of reads, a genome containing all genetic information about an organism will be constructed. In the genome assembly, reads(substrings) are concatenated together to construct the genome. The next generation sequencing (NGS), which is one of the sequencing technology, output (reads) genome reads are available in public repositories such as National Centre for Biotechnology Information (NCBI) [54], and DNA Databank of Japan (DDBJ), UC Santa Cruz Genomics Institute (UCSC) [55]. Genome data reads can be stored in FASTA and FASTQ files [54].

Data Structure Approaches for Genome Assembly

An exact match is considered when two matched strings are identical, as shown in Figure.2-A. On the other hand, the approximate match represents the high similarity between two strings and is not necessarily identical, and some differences might be there. The number of different letters between similar approximate matched strings is called distance [5]. The distance can be a substitution representing a letter in the target string differently from its corresponding one in the reference string [43]. The distance can also be an insertion of a letter in the target string that does not exist in its corresponding offset in the reference string. Another type of distance is deletion, which one letter is absent in the target string compared to its corresponding offset in the reference string, as shown in Figure.2-B. *De novo* is Latin word that means from scratch or new. One of the genome assembly approaches is de novo approach, which concatenates overlapping reads with the exact match between the suffix of a read with a prefix of another read [6]. Overlapping reads will be concatenated into a longer string called a contig. Similarly, overlapping contigs will be concatenated into a scaffold, as shown in Figure 3.[7]. Overlapping distance (length) is a significant parameter for any overlapping algorithm. Overlapping errors have a small impact when sufficient distance is used, as shown in Figure.4, which suggests that empirical overlapping distance should be greater than 40% from the read length [8]. On the other hand, the assembly is guided by known reference genome instances in the reference-guided approach. Reads are mapped against the reference genome to determine their orders, called the reads alignment process [46], as shown in Figure.5. The reads alignment against the reference genome is based on an approximate match between the read and its corresponding k-mer of the reference genome. Mer is a Latin word that means part, while k is the length of this part [9]. It is not necessarily the reads exactly match the reference k-mer, because reads of the target genome are slightly different from the used reference genome due to individual biological variation [1].

Limitations Of Current Sequencing Technology as the Root Cause of Genome Assembly Challenges

In 1965 Fred Sanger and colleagues sequenced the first DNA reads to produce reads slightly less than one kilo-base (kb) in length. However, it was first-generation DNA sequencing. Second-generation DNA sequencing or next-generation sequencing (NGS) is a massive number of tiny Sanger sequencing running in parallel. In NGS, large quantities of DNA can be sequenced quickly. Short-length reads are considered the first limitation in this technology, generating a big challenge for assembling repetitive sequences from short reads. Recently, third-generation DNA

sequencing (TGS) was introduced by Stephen Quake. The main characteristic of this method is generating longer-length reads. TGS is currently available, generating longer-length reads able to bridge long repetitive sequences, and used to develop assemblies such as Canu and Flye [65]. However, the natural high error rate of up to 15% of long reads impacts the accuracy of the assembly requiring higher computational complexity correction processes when processing large datasets [57],[64]. The most common question is why DNA molecules cannot be sequenced as a single string. Firstly, the enzyme used in the sequencing experiment naturally generates sequence errors after sequencing three kilo-base [11]. Secondly, the DNA sequencing experiment is too much time-consuming. Thus, sequencing the entire genome using a single sequencing machine would take years to be done [12]. There is no sequencing platform that provides complete coverage of the whole genome as a single string [13]. In NGS, each region is sequenced many times in order to ensure that genome is accurately sequenced, called coverage. Data coverage is the average number of times a base of a genome is sequenced to ensure that the genome is sequenced accurately. It is often expressed as (1x, 2x, 3x,..., nx) as shown in Figure.6. Target genome reads from NGS are extremely massive data. For example, the number of 60X sequence coverage is often mentioned in that the number of letters in the dataset is 60 times bigger than the original genome size [14]. This massive data amount is considered the second limitation of NGS, resulting in big performance challenges during genome assembly.

Genome Assembly Challenges

Assembling Repeats from Short Reads Misassembly

The presence of repeats in a genome has increased the complexity of accurately assembling the genome from short reads. [15],[35]. In the *de novo* approach based on overlapping, algorithms might concatenate reads wrongly, leading to a misassembly simply because they have the same matched overlapping distance [9],[34]. The formulation in Figure.7 assumes that the given genome $G=\{TGGGACTGG\}$, then the genome G is sequenced, resulting in an array of reads $R=\{GACT, ACTGG, TGGGAC\}$. According to computer science, concatenating these reads based on overlapping is correct if and only if the reads have an exact suffix-prefix match at a specific overlapping distance. Thus, so both assembled genomes $\{GACTGGGAC\}$ and $\{TGGGACTGG\}$ are correct solutions accordingly. However, because of the ambiguity that is resulted from repeats in short reads, when the genome is assembled genome to be $\{GACTGGGAC\}$ while the real one is $\{TGGGACTGG\}$, this problem is a well-known problem called misassembly [17].

Most *de novo* overlapping algorithms have a high degree of misassembly resulting from repeat ambiguity. The greedy algorithm is one of the simplest algorithms in a data structure. It is used for overlapping. It starts with joining the two reads that overlap the highest. Unassembled reads are chosen for the next run until a pre-defined minimum threshold is reached [9]. Although the greedy approach is computationally feasible, it cannot distinguish the repeats, merging reads

wrongly resulting in misassembly [9] as shown in the example in Figure.8 for the Genome $G=\{GGATGGGGATGCCT\}$ that has two repeats "ATG". Overlap Layout Consensus (OLC) is a famous graph-based algorithm introduced by Staden in 1980 and subsequently improved and extended by others. Many *de novo* assembly tools are based on OLC methods, such as Newbler, PCAP, Celera Celera, CAP3, and ARACHNE. In OLC, repeats in short reads create a mathematical problem of finding a Hamiltonian path. The graph will be too complex if the genome is highly repetitive, such as the human genome.

Another widely used algorithm is De Bruijn Graphs (DBG), which Nicholas Govert de Bruijn introduced in the 1940s. DBG considers each read is a k-mer from the Genome [9]. The first step in DBG is to divide each read with length k into two substrings (edges) with length (k-1). Next, draw each left substring corresponding to its right substring. A path that contains every edge of the graph exactly once is called the Eulerian path or Eulerian walk. DBG is a widely used algorithm. Current de novo assemblers use DBG, such as Euler, Velvet, SOA Pdenovo, ABySS, and IDBA. However, repeats in short reads create a fragmented assembly problem in this algorithm, shown in Figure.9.

In the reference-guided approach, repetitive sequences can result in ambiguity during the alignment process when a read has approximately matches with many positions (offsets) in the reference genome. This ambiguity is a well-known data structure problem called the multiple-matching problem. However, multiple-matching results in misassembly, as shown in Figure.10. Most alignment tools cannot accurately map reads against the reference genome. The sensitivity or confidence of the Alignment tool (S) is calculated by the ratio of correctly mapped read to incorrectly mapped read at a particular threshold ($S = \text{Number of reads mapped correctly} / \text{a number of reads mapped incorrectly}$). Repeats in short reads affect the sensitivity of the alignment tool, resulting in a degree of misassembly, decreasing the total accuracy of assembly [18].

Computational Performance Challenges in Genome Assembly

The second computational problem in genome assembly is challenging performance due to the massive data amount generated by high coverage NGS. This massive data amount generates big performance challenges in each assembly approach in terms of CPU, memory, and storage. Assembly performance challenges exist in both assembly approaches [14]. In the reference-guided approach, reads alignment works in time complexity $O(NL)$. N is the number of reads (millions) to be aligned against each offset of reference genome with length L, which might be billions of letters. It is a highly intensive process, as shown in Figure.11. On the other hand, computational complexity in *de novo* overlapping works in time complexity $O(N)^2$, where (N) is the number of reads (millions). Comparing millions of reads against millions to detect overlapping is a well-known data structure querying problem called all-against-all or all-suffix-prefix problem (ASPP), as shown in Figure.11 [8].

Accuracy And Performance Evaluation in Genome Assembly

The accuracy of the genome assembly produced through the assembly approaches can be evaluated using metrics scores throughout the assembly process, such as counting, the number of contigs, the proportion of reads that can be aligned against the known reference genome if exists (misassembly degree), and the absolute length of contigs. A commonly used statistical metric is N50. Given a set of contigs, each with its own length, the N50 is defined as the shortest contigs' length in the group of contigs which represent 50% of the assembly length. N50 can be used to denote the contiguity of the assembly. The pseudocode of how to calculate N50 from an array of contigs is shown in the pseudocode in Box.1 [48]. Although N50 is still widely used for assembly evaluation, this metric does not reflect the quality of an assembly and can be misleading [21]. An example of using N50 to evaluate the genome assembly is a *Ciona intestinalis* genome with an estimated N50 of 234 kb length, which was found inaccurate a few years later due to repeat problems that resulted in misassembly. The N50 metric does not consider that some contigs may be erroneously joined or even overlapped [58] [62]. U50 is similar to N50 in the calculation, and the only difference is that U50 uses unique contigs for calculation that do not overlap with other contigs to represent a more accurate result [21]. Another method to ensure the completeness of an assembly is by detecting the presence-absence variations (PAV) against the reference genome. It identifies the sequences in the reference genome that are entirely missing in the newly generated assembly [58]. Also, QUAST-LG is a tool that compares large genomic de novo assemblies against reference sequences and calculates the quality metrics. [59]. The Merqury, is a reference-free tool assembly evaluation By comparing k-mers in a *de novo* assembly to those found in unassembled high-accuracy reads [60][61].

The performance of genome assembly can also be evaluated through the computational execution time (time complexity), referred to as time complexity, which is represented as $O(N)$, where (N) is the size of the input and (O) is the execution time for the algorithm, it can be linear, logarithmic, quadratic, exponential or any other relationship. For example, if the number of reads to search for overlapping is five reads, and the index hits ten times, the time complexity of overlapping is $O(2N)$, as shown in the pseudocode in Box.2. Also, the computational performance of the genome assembly can also be assessed by memory or storage consumption.

Survey methodology

This review article summarizes the enhancement in genome assembly accuracy using prior repeat identification and hybrid assembly approaches. Also, we highlight the limitations of current methods. On the other hand, we overview the most famous data structure indexing to optimize genome assembly, such as k-mer index, tree index, hash index, also parallelization methods. We also highlight the limitations and show that combining more than one method improves accuracy and performance. This review helps computer science researchers understand genome data and its considerations, such as repeat. Also, it will help to identify the main research problems in the genome assembly area in terms of accuracy and performance optimization.

I. Search Strategy

The authors found the primary studies through seven digital libraries: Web of Science, Science Direct, SCOPUS, Google Scholar, IEEE Explore, PubMed, and EMBASE. These digital libraries are most preferred and widely used by the research community. The search query string is applied to an individual database to find relevant studies based on the research questions shown below in box.3. The search strategy for the literature process is shown in Figure.12. A total of 65 research papers, after the duplications were removed, were found and downloaded from the mentioned database.

II. Inclusion criteria

The studies concerning that are:

- A. Written in English
- B. Based on genome assembly.
- C. Focused on accuracy and performance optimization.
- D. Published in conference or journal

III. Current Solutions for Misassembly Problem

A. Prior Repeat Identification Methods

Repeat identification is a method to detect and identify repetitive sequences in the reads or in the reference genome [1]. Identifying repetitive sequences before genome assembly helps the assembly algorithm distinguish the reads that have repeats, and avoid misassembly, thus, significantly enhancing the accuracy of genome assembly from short reads [23]. Repeat identification methods are introduced first in the reference-guided approach through the repeat masker method. However, it did not work well with big genomes [1]. The de novo approach is preferable, identifying repeats from target genome reads rather than a reference genome [6]. Most of the current repeat identification methods are statistically based, calculating the occurrences of repetitive sequences in reads. In order to clearly understand how it works, let (G) is the genome that is sequenced with 5X coverage, so ($C=5$). Let rs is a repetitive sequence with a length of 9 letters, $rs=\{ACGTGATAT\}$. Let R denote the array of reads generated by sequencing genome (G), $R=\{r_1, r_2 \dots r_i\}$ as shown in Figure.13. Reads are just a substring of Genome (G). According to the current sequencing coverage, we expect any unique substring of (G) to appear in (R) several times less than or equal to (C). In this example, a read, $r_i\{GTG\}$ appears in the data set six times, which means the frequency of r_i , $f_{ri}=6$, while the sequencing coverage $C=5$. If r_i sequence is a unique sequence in Genome (G) the $\max f_{ri}=5$ under coverage ($C=5$) and $\max f_{ri}=6$ under coverage ($C=6$), and so on. If ($f_{rs}=1$) in (G), it is impossible that the frequency of the entire rs or even any substring of rs exceeds the sequencing coverage (C) in the reads dataset. In Figure. 13, $f_{ri}=6$, greater than C, simply means that the extra occurrence of r_i comes from sequencing another copy of rs. If the read frequency is less than or equal to the sequencing coverage, it means that this read comes from a unique sequence in the genome. On the other hand, if the frequency of the read exceeds the well-known sequencing coverage, it

means that this read comes from the repetitive sequence in the genome. More than half of some genomes are repeats and vary in length, which might be too long, such as TE. Thus the size of repetitive sequences is used as a criterion to evaluate the accuracy of identifying repeats in many repeat identification methods [3][63][64]. Many algorithms have been proposed according to the previous statistical concept. Earlier, the RePS algorithm assumed that any 20-mer (substring from the read with length 20 letters) that appears in the dataset more often than a sequencing coverage is likely to be an exact repeat and is masked out. A sliding window assembly (SWA) is a genome assembly algorithm proposed by [24] for identifying and assembling repeats. Fundamentally, the algorithm considers the entire read is a substring of the genome. SWA calculates the frequency of each read in the dataset. Based on the well-known sequencing coverage SWA splits reads into two groups, repeat and non-repeat, assembling them separately. The non-repeat group consists of those reads with a frequency less than the sequencing coverage, while the repeat group includes those with a frequency greater than the sequencing coverage. Then, the overlap is run for each group separately in parallel to construct the genome superstring. However, SWA considers the entire reads as the substring of genome superstring, which means it can detect repeat length that is only greater than or equal to the read length, while repetitive sequences might be smaller than the read's length, such as microsatellite tandem repeat, which cannot be detected by SWA [1] as shown in Figure.14. Repeat *De novo* (REPdenovo) is a repeat identification method proposed by [25], similar to the earlier method RepARK that splits reads into smaller substrings (k-mers) then identifies the frequency of k-mers, group them to repeat and non-repeat groups, finally, assembles these k-mers. Unlike SWA, REPdenovo uses the average k-mers frequency as a threshold of this algorithm, where a repeat is identified when a k-mer frequency is more than the threshold, as shown in Figure.15. REPdenovo is also able to identify a repeat shorter than the read length by constructing them from reads' k-mers to longer repeats. Evaluation and comparison show that REPdenovo outperforms the earlier RepARK method regarding the accuracy and completeness of repeats construction. REPdenovo discovered that previous repeat annotations had missed a significant number of 190 potentially new repeats in the human genome. The high accuracy of identifying the repeat enhances the genome's contiguity with N50 3141, while RepARK method only N50 116. However, when reads are chopped up into k-mers the biological info represented in the reads is lost [12]. Similar to REPdenovo, the Detection of Long Repeats (DLR) converts all reads into unique k-mers of a certain length and screens out the k-mers with a high frequency [26]. The Detection method in DACCOR assembly is also similar to REPdenovo calculation introduced as a stage in Characterization Reconstruction (DACCOR) hybrid assembly. Unlike REPdenovo, it splits the reference genome into k-mers instead of the reads themselves, then identifies repeats from k-mers of the reference genome, similarly to REPdenovo calculation. The detection method in DACCOR assembly identifies repeats from the reference genome, which might come with low identification accuracy, simply because repeats in the reference genome might be slightly different from the real repeat in target genome reads [1].

On the other hand, digital signal processing can also detect repetitive elements. In most signal-processing methods, DNA sequences are converted to numerical sequences, and repetitive

elements can be identified by Fourier power [49]. However, capturing the essential features of repetitive elements, such as copy numbers of repeats, is still challenging. In addition, Fourier transform cannot capture repeats in short genomes [49].

Recently, machine learning **has been employed for** identifying repetitive elements. Repeat Detector (Red) has been proposed by [27] as de novo tool for discovering repetitive elements in genome reads. Red utilizes a Hidden Markov model (HMM) dependent on labeled training data, and it successfully identified new repeats in the Human Genome. However, it only generates genome coordinates for repeats without any annotation. Therefore, the red output does not help analyze repeat content or TEs evolution[1]. **Improvements are required in optimizing computation requirements and choosing suitable training datasets for such a machine-learning program in order to strengthen the program accuracy and precision on the repeat prediction in the genomes** [28].

As shown in the current limitations of some current repeat identification methods in Table.1. No one algorithm fits all lengths of repetitive sequences [29]. The combination of two or more methods might introduce a better solution to identify and assemble repeat accurately. A combination of repeat identification methods consistently outperforms the usage of a single method due to the variety of repetitive sequences. An example of this combination is merging consensus sequences generated by RepeatModeler and Repbase library, which successfully annotated many repeats in many genomes [1]. Another example of a repeat identification methods combination is Tandem Repeats Finder (TRFi), used to identify tandem repeats and combined with the RepeatExplorer method, and the result outperforms both of them in identifying complex tandem repetitive sequences [30]. However, combining methods adds more computational complexity, which might be addressed by data indexing and increasing the parallelism level during the identification process.

B. Hybrid Assembly Approach

The last two decades have seen significant improvement in solving the misassembly problem that resulted from repeats through the hybridization concept. Fundamentally, hybridization combines two or more methods to perform genome assembly. Hybridization appears in genome assembly in different ways. It can combine data from different platforms [44] or two or more genome assembly methods. The example is combining the OLC and the DBG method to introduce DBG2OLC to enhance assembly accuracy that is affected by repeats [13]. The hybridization between de novo and reference-guided approaches creates a hybrid assembly hybrid approach which able to complement each approach's limitation [32]. The idea of a hybrid assembly approach has been introduced by [33], who proposed a scaffold-builder assembler composing an initial de novo assembly from reads , followed by aligning contigs against a reference genome to build a scaffold. Enhancing assembly accuracy through the hybrid assembly approach adds many contributions to understand the diversity of species.. An example of is a study on Cyclospora cayetanensis species, , a parasite that caused intestinal infection, has a highly repetitive genome that could not be assembled accurately using de novo alone. Combining a reference-guided approach with a de novo approach generates new assembly for this species with a significantly greater depth of coverage and a lower degree of misassembly [7].

Although, accurate reconstruction of repetitive sequences cannot be without repeat resolution (repeat identification). Most of the current assembly hybrid approaches do not employ prior repeat identification [23]. The accuracy is enhanced by using more than one reference genomes or multiple samples of the target genome from different sequencing platforms [7]. The characterization and reconstruction of repetitive regions (DACCOR) introduced a new idea to enhance the hybrid assembly approach based on identifying repetitive sequences in the reference genome prior to hybrid assembly. Although the detection repeat identification method in DACCOR identifies repetitive sequences from the reference genome, the results on the bacterial genome *Treponema pallidum*, as shown in Table.2, reveals high accuracy in identifying repetitive sequences [23]. In Table.3, the comparison of the accuracy in generating contigs from repetitive regions detecting repeats for *T. pallidum* (sample AR1 with the coverage X157), DACCOR successfully constructed 96.8% of contigs that can be mapped to the reference genome, while, SPAdes de novo assembly generates about 69.5% of contigs from repetitive regions that can be mapped to the reference genome and EAGER reference guided approach with generates around 42.1%. Using repetitive sequences from the reference genome to detect repeats in the target genome in DACCOR might have advantages over the conventional assembly approach. However, the repeats identified through the approach might be limited to the reference genome, which may not well represent the actual repetitive sequences present in the target genome. On the other hand, the hybrid assembly approach comes with big performance challenges. The total assembly time in the hybrid approach is exceptionally challenging for any traditional algorithm or hardware architecture. Although performance challenges are well-optimized in each approach individually [36], no data structure indexing is introduced to fit overlapping and reads alignment together in the hybrid assembly approach.

IV. Current Solution For Performance Optimization Of Genome Data

A. Indexing in Genome assembly

Indexing could speed up the search by reducing the number of iterations (time complexity), leading the algorithm to focus on a specific part of the data instead of searching in the entire dataset [50]. One of the earliest data structure indexing methods is k-mer index. However, massive genome data requires more optimization for this index. The keys of k-mer index can be extracted as a shape, skipping letter instead of extracting all k-mers from the string to speed up the index hitting. As shown in Figure.16, k-mers are extracted as a pattern by taking the first and third letters followed by the fifth and sixth letters. This reduces the index total size, which improves the specificity of the index hitting and index verification. This kind of index must be hit with the same index pattern. This idea was implemented in Homopolymer Compressed k-mers method proposed by [36]. This implementation added significant enhancement of mapping reads to the reference genome. Another variation to speed up querying k-mer index is a binary search. Binary search, also known as half-interval search, logarithmic search is a search algorithm that finds the position of a target value within a sorted array. The k-mers of the genome must be ordered alphabetically in ascending order, as shown in Figure.17. In this case, the query does not look up in the entire index simply because (TGG) is alphabetically greater than (GTG). The query does bisection (dividing the problem into two halves), skipping the first part of the index [31]. Each iteration bisection occurs, resulting in a significant reduce the query time. The total number of bisections that are needed to perform can be calculated as $\text{Log}_2(n)$.

This implementation significantly reduces search trials during reads mapping to reference genome and can be implemented for reads during the search for overlapping [51].

The Hash index is the most famous implementation of k-mer index. In the hash index, the hash function is responsible for mapping each distinct k-mer to one bucket (group), as shown in Figure.18. the Hash index is widely used in many studies on genome data for indexing dataset reads [10] to reduce suffix-prefix search in de novo overlapping. However, because of repeat, sometimes two keys may generate an identical hash causing both keys to point to the same bucket, which is known as a hash collision, slowing down the search in the hash index [50].

Suffix indexes are another data structure index family. There are four variations of suffix indexes: trie suffix tree, suffix array, and FM index. However, the suffix tree (ST) is the most popular and widely used indexing tool in bioinformatics applications [16]. It is the base of some popular sequence alignment tools, such MUMmer and RePuter, and can be used to solve more complex problems such as repeat identification. Unfortunately, the construction of ST is highly memory and CPU-consuming [52]. However, there are many research efforts to reduce the construction of ST time complexity, for example, scaling the construction of ST on multiple CPU cores (64 cores) that [37] achieves a speedup from 2X to 4X over the original algorithm to construct S.T. Figure.19 shows the construction of ST index from genome $G=\{ACGCGT\}$.

In overlapping massive genome reads, maybe performing an approximate match between suffix and prefix first add a significant reduction of time complexity of overlapping by skipping the exact match for suffix prefix candidates that do not have an approximate match. The pigeonhole principle (PP) is a well-known data structure principle called seed and extends principle that can achieve this idea. PP separates the read into non-overlapping partitions, and the exact matching algorithm checks the match of one of the partitions considered suffix prefix approximate match, then the verification step must be done as shown in Figure.20. According to this powerful principle, many algorithms and methods for indexing introduced to genome data took advantage of the pigeonhole principle and suffix-tree index. The result of *Escherichia Coli* genome shows that the pigeonhole solution with prefix tree is superior in terms of time and storage compared to the traditional suffix tree for overlapping [8], as shown in Table.4.

The optimization using data structure indexes for overlapping reads and read alignment in the genome hybrid assembly approaches remains a challenge. However, ST is the best solution for repetitive sequences scenarios, its construction memory, and CPU consumption. In Hash, index collisions slow down the search at this key. FM.- index is based on Burrows-Wheeler Transformation (BWT) that was proposed by Burrows and Wheeler in 1994 [38]. Querying FM index cannot use binary search because it does not consist of the complete rotation of the genome and needs another process to find the offset of the index hit, which is a highly computational task in massive genome data.

Hybridization between two indexes might introduce a better solution for indexing massive genome data. An example of the hybrid between two indexes is MiniSR index which uses k-mer indexing with the hash index. MiniSR can index the massive human genome in a few minutes [53]. The hybridization between a few indexed might introduce a solution to optimize overlapping and read alignment in hybrid genome assembly.

B. Parallelization in Genome Assembly

The exponential growth of genomic data has instigated a significant challenge for genomics analysis computing infrastructure and software algorithms. Genomic experiments are now reaching the size of Terabytes and Petabytes [22]. Scientists may require weeks or months to process this massive amount of data using their own workstations [22]. Parallelism techniques and high-performance computing (HPC) environments can help to reduce the total processing time [20], [22]. The simplest definition of parallelization is to split the extensive process into smaller subprocesses and run them in parallel on multiple CPUs [52][22]. It solved many performance challenges in the big data area. However, parallelization is introduced in genome data for genome assembly in some stages individually.

Sequencing technologies produce millions to billions of short reads. The first step to assembling these reads is to extract them from FASTA or FASTQ files into structured database tables, which is too time-consuming. Parallelization could reduce raw data streaming time is to split big FASTA or FASTAQ files into smaller files processed in parallel [56].

Another variation of parallelization in genome data is parallelizing the index's construction. An example of this idea is proposed by [39] for speeding up prefix tree index construction by letting each processor work on strings that start with a specific character in the alphabet (A, C, G, and T). For example, processor 1 constructs the part of the tree that corresponds to strings that start with the letter "A" while processors 2, 3, and 4 construct the parts of the tree corresponding to the strings starting with "C", "G", and "T" respectively. A similar idea was introduced by [40], proposing HipMer assembler to construct k-mer index using a deterministic function to map each k-mer to a target processor.

In de novo overlapping and read alignment of reference-guided approach, the IBD algorithm addressed pairwise comparison problems that can be used to overlap two reads or align reads against k-mer of the reference genome [19]. The algorithm breaks the problem (N^2) using massive parallelization with each order (N) comparison. This approach is applied to the U.K. Biobank dataset, with a 250X faster time and 750X less memory usage over the standard approach of pairwise alignment.

Over the last decades, many computing platforms parallelize different stages of genome assembly. For instance, Hadoop[22] introduced a powerful idea of parallelizing the suffix tree index construction. However, the Apache spark platform works 100 times faster than Hadoop, especially in iterative operators [41]. A distributed and parallel computing tool named HAlign-II was introduced by [22] to address read alignment with extremely high memory efficiency.

GraphSeq method proposed splitting genome files and using a high scaling platform [42]. GraphSeq works on the Apache spark platform and achieved 13X speedup *de novo* genome assembly. GraphSeq splits a big compressed file into several small ones loading all of the reads in parallel, generating the corresponding suffixes, grouping those suffixes with the exact initial string into the same partition, and applying string graph construction in parallel by partitions. Similarly, an innovative algorithmic approach proposed by [2] is called Scalable Overlap-graph Reduction Algorithms (SORA), performing string graph reduction on Apache Spark. SORA was evaluated with human genome samples to process a nearly one billion edge of string graph in a short time frame with linear scaling.

Parallelization enhances the performance of streaming massive genome data, genome index construction, and overlapping in the *de novo* approach or read alignment in the reference-guided approach. Despite its performance, due to the nature of the hybrid assembly approach, to our best knowledge, the usage of parallelization in overlapping and reads alignment for genome assembly is yet to be fully implemented.

Conclusions

Assembling genome reads with the presence of repetitive sequences at a good quality is essentially important but, indeed, far more challenging. Ignoring repeats in the genome would revoke the purpose of the genome assembly. A massive number of reads produced by sequencing technologies have caused high computational performance for the genome assembly approach. A lot of research has contributed significantly to the development of enhancing assembly accuracy by reducing the degree of misassembly in the genome assembly approaches. In this paper, we suggest by using prior repeat identification methods to create prior repeat resolution guiding overlapping algorithm in *de novo* or alignment algorithm in the reference-guided approach could reduce the chance of genome misassembly. However, none of the repeat identification methods can accurately detect all types with different lengths of repetitive sequences. The combination of *de novo* with the reference-guided assembly approach could yield better results for genome assembly, which can be another alternative solution. Although the hybrid assembly approaches found to outperform the individual *de novo* approach and the reference-guided approach, repetitive sequences can still generate misassembly, which could be enhanced through the use of

multiple samples or multiple reference genomes. However, using multiple samples or reference genomes will add high computational complexity to the hybrid assembly approach. Moreover, most of the current assembly hybrid approaches still lack repeat identification prior to the approach, which can help to increase the accuracy of the genome assembled.

The combination between the *de novo* approach and the reference-guide approach comes with a more computationally extensive performance challenge than using the approach individually. Therefore, the investigation of optimizing the computational performance challenges in the hybrid assembly approach for overlapping and reads alignment is still an open challenge.

Enhancement of the repeat identification accuracy by approaches would be one of the main priorities in addressing issues of genome misassembly. The enhanced repeat identification method should be added at the prior stage for hybrid genome assembly to enhance the accuracy without relying on multiple samples or reference genomes, which can increase the computational complexity of the the hybrid assembly approach.

Consequently, the performance of the hybrid assembly approach might be optimized through the hybridization of indexing methods with parallelization to optimize overlapping and reads alignment in the hybrid genome assembly approach.

Hybrid genome assembly enhanced by accurate prior repeat identification and optimized by a combination of indexing and parallelization, as shown in Figure.21, might be a novel solution for assembling repetitive genomes from short reads.

Conflicts of Interest

The authors declare no conflict of interest.

ACKNOWLEDGMENT:

This research was supported by the Ministry of Higher Education (MOHE) through Fundamental Research Grant Scheme FRGS/1/2019/ICT02/UTM/02/13, Research anagement Centre (RMC), UTM, and ALI@S research group.

References

1. Zeng, L., Kortschak, R.D., Raison, J.M., Bertozzi, T. and Adelson, D.L., 2018. Superior ab initio identification, annotation and characterisation of T.E.s and segmental duplications from genome assemblies. PloS one, 13(3), p.e0193588.
2. Paulson, H., 2018. Repeat expansion diseases. Handbook of clinical neurology, 147, pp.105-123.
3. Genovese LM, Geraci F, Corrado L, Mangano E, D'Aurizio R, Bordoni R, Servergnini M, Manzini G, Bellis GDe, D'Alfonso S, Pellegrini M. 2018 "A census of Tandemly repeated polymorphic loci in genic regions through the comparative integration of human genome

- assemblies,” *Frontiers in Genetics*, 9, p.155. Available at:
<https://doi.org/10.3389/fgene.2018.00155>.
4. Pinto, B.S., Saxena, T., Oliveira, R., Méndez-Gómez, H.R., Cleary, J.D., Denes, L.T., McConnell, O., Arboleda, J., Xia, G., Swanson, M.S. and Wang, E.T., 2017. Impeding transcription of expanded microsatellite repeats by deactivated Cas9. *Molecular cell*, 68(3), pp.479-490.
5. Patil, N., Toshniwal, D. and Garg, K., 2013. Genome data classification based on fuzzy matching. *CSI Transactions on ICT*, 1(1), pp.9-28.
6. Baichoo, S. and Ouzounis, C.A., 2017. Computational complexity of algorithms for sequence comparison, short-read assembly and genome alignment. *Biosystems*, 156, pp.72-85.
7. Gopinath, G.R., Cinar, H.N., Murphy, H.R., Durigan, M., Almeria, M., Tall, B.D. and DaSilva, A.J., 2018. A hybrid reference-guided de novo assembly approach for generating *Cyclospora* mitochondrion genomes. *Gut Pathogens*, 10(1), pp.1-8. Available at:
<https://doi.org/10.1186/s13099-018-0242-0>
8. Haj Rachid, M. (2017) “Two efficient techniques to find approximate overlaps between sequences,” *BioMed Research International*, 2017, pp. 1–8. Available at:
<https://doi.org/10.1155/2017/2731385>.
9. Simpson, J.T. and Pop, M., 2015. The theory and practice of genome sequence assembly. *Annual review of genomics and human genetics*, 16, pp.153-172.
10. Zhang, H., Chan, Y., Fan, K., Schmidt, B., & Liu, W. (2018). Fast and efficient short read mapping based on a succinct hash index. *BMC bioinformatics*, 19(1), 1-14.
11. Garibyan, L. and Avashia, N., 2013. Research techniques made simple: polymerase chain reaction (PCR). *The Journal of investigative dermatology*, 133(3), p.e6.
12. Angeleska, A., Kleessen, S. and Nikoloski, Z., 2014. The Sequence Reconstruction Problem. In *Discrete and Topological Models in Molecular Biology* (pp. 23-43). Springer, Berlin, Heidelberg.
13. Jain, M., Olsen, H.E., Turner, D.J., Stoddart, D., Bulazel, K.V., Paten, B., Haussler, D., Willard, H.F., Akeson, M. and Miga, K.H., 2018. Linear assembly of a human centromere on the Y chromosome. *Nature biotechnology*, 36(4), pp.321-323. Available at:
<http://dx.doi.org/10.1038/nbt.4109>
14. Jain, M., Olsen, H.E., Turner, D.J., Stoddart, D., Bulazel, K.V., Paten, B., Haussler, D., Willard, H.F., Akeson, M. and Miga, K.H., 2018. Linear assembly of a human centromere on the Y chromosome. *Nature biotechnology*, 36(4), pp.321-323. Available at:
<http://dx.doi.org/10.1038/nbt.4109>
15. Lohmann, K. and Klein, C., 2014. Next generation sequencing and the future of genetic diagnosis. *Neurotherapeutics*, 11(4), pp.699-707.
16. Barsky, M., Stege, U., Thomo, A., & Upton, C. (2009, November). Suffix trees for very large genomic sequences. In *Proceedings of the 18th ACM conference on information and knowledge management* (pp. 1417-1420).
17. Medvedev, P., 2019. Modeling biological problems in computer science: a case study in 5genome assembly. *Briefings in bioinformatics*, 20(4), pp.1376-1383.

18. Thankaswamy-Kosalai, S., Sen, P. and Nookaew, I., 2017. Evaluation and assessment of read-mapping by multiple next-generation sequencing aligners based on genome-wide characteristics. *Genomics*, 109(3), pp.186-191.
19. Sapin, E., & Keller, M. C. (2021). Novel approach for parallelizing pairwise comparison problems as applied to detecting segments identical by decent in whole-genome data. *Bioinformatics*, 37(15), 2121-2125.
20. Ocaña, K. and de Oliveira, D., 2015. Parallel computing in genomic research: advances and applications. *Advances and applications in bioinformatics and chemistry*, pp.23-35. Available at: <https://doi.org/10.2147/AABC.S64482>
21. Lischer, H.E. and Shimizu, K.K., 2017. Reference-guided de novo assembly approach improves genome reconstruction for related species. *BMC bioinformatics*, 18(1), p.474.
22. Shi, L. and Wang, Z., 2019. Computational strategies for scalable genomics analysis. *Genes*, 10(12), p.1017. Available at: <https://doi.org/10.3390/genes10121017>
23. Seitz, A., Hanssen, F. and Nieselt, K., 2018. DACCOR—Detection, characterization, and reconstruction of repetitive regions in bacterial genomes. *PeerJ*, 6, p.e4742.
24. Lian, S., Li, Q., Dai, Z., Xiang, Q. and Dai, X., (2014) “Ade novogenome assembly algorithm for repeats and nonrepeats,” *BioMed Research International*, 2014, pp. 1–16. Available at: <https://doi.org/10.1155/2014/736473>.
25. Chu, C., Nielsen, R. and Wu, Y., 2016. REPdenovo: inferring de novo repeat motifs from short sequence reads. *PloS one*, 11(3), p.e0150719.
26. Liao, X., Zhang, X., Wu, F.X. and Wang, J., 2019, November. de novo repeat detection based on the third generation sequencing reads. In 2019 IEEE International Conference on Bioinformatics and Biomedicine (BIBM) (pp. 431-436). IEEE.
27. Girgis, H.Z., 2015. Red: an intelligent, rapid, accurate tool for detecting repeats de-novo on the genomic scale. *BMC bioinformatics*, 16(1), p.227.
28. Libbrecht, M.W. and Noble, W.S., 2015. Machine learning applications in genetics and genomics. *Nature Reviews Genetics*, 16(6), pp.321-332.
29. Guo, R., Li, Y.R., He, S., Ou-Yang, L., Sun, Y. and Zhu, Z., 2018. RepLong: de novo repeat identification using long read sequencing data. *Bioinformatics*, 34(7), pp.1099-1107. Available at: <https://doi.org/10.1093/bioinformatics/btx717>
30. Peška, V., Sitová, Z., Fajkus, P. and Fajkus, J., 2017. BAL31-NGS approach for identification of telomeres de novo in large genomes. *Methods*, 114, pp.16-27.
31. Berztiss, A. (2014) *Data structures: Theory and practice (Computer science and applied mathematics)*. Ex-library. Academic Press. New York, USA.
32. Platt, R.N., Blanco-Berdugo, L. and Ray, D.A., 2016. Accurate transposable element annotation is vital when analyzing new genome assemblies. *Genome biology and evolution*, 8(2), pp.403-410.
33. Silva, G.G., Dutilh, B.E., Matthews, T.D., Elkins, K., Schmieder, R., Dinsdale, E.A. and Edwards, R.A., 2013. Combining de novo and reference-guided assembly with scaffold_builder. *Source code for biology and medicine*, 8(1), p.23.
34. Wang, P., Meng, F., Moore, B. M., & Shiu, S. H. (2021). Impact of short-read sequencing on the misassembly of a plant genome. *BMC genomics*, 22(1), 1-18.

35. Acuña-Amador, L., Primot, A., Cadieu, E., Roulet, A., & Barloy-Hubler, F. (2018). Genomic repeats, misassembly and reannotation: a case study with long-read resequencing of *Porphyromonas gingivalis* reference strains. *BMC genomics*, 19(1), 1-24.
36. Liu, Y., Yu, Z., Dinger, M.E. and Li, J., 2018. Index suffix-prefix overlaps by (w, k)-minimizer to generate long contigs for reads compression. *Bioinformatics*, 35(12), pp.2066-2074.
37. Labeit, J., Shun, J. and Bletloch, G.E., 2017. Parallel lightweight wavelet tree, suffix array and FM-index construction. *Journal of Discrete Algorithms*, 43, pp.2-17.
38. Wang, J., Wong, G.K.S., Ni, P., Han, Y., Huang, X., Zhang, J., Ye, C., Zhang, Y., Hu, J., Zhang, K. and Xu, X., 2002. RePS: a sequence assembler that masks exact repeats identified from the shotgun data. *Genome research*, 12(5), pp.824-831.
39. Haj Rachid, M. and Malluhi, Q. (2015) "A practical and scalable tool to find overlaps between sequences," *BioMed Research International*, 2015, pp. 1–12. Available at: <https://doi.org/10.1155/2015/905261>.
40. Ellis, M., Georganas, E., Egan, R., Hofmeyr, S., Buluç, A., Cook, B., Olikier, L. and Yelick, K., 2017, August. Performance Characterization of De Novo Genome Assembly on Leading Parallel Systems. In *European Conference on Parallel Processing* (pp. 79-91). Springer, Cham.
41. Wan, S. and Zou, Q., 2017. HAlign-II: efficient ultra-large multiple sequence alignment and phylogenetic tree reconstruction with distributed and parallel computing. *Algorithms for Molecular Biology*, 12(1), pp.1-10. available at: <https://doi.org/10.1186/s13015-017-0116-x>
42. Su, C.T., Chang, M.T., Cheng, Y.C., Li, Y.L. and Wang, Y.T., 2018. GraphSeq: Accelerating String Graph Construction for De Novo Assembly on Spark. *bioRxiv*, p.321729." Available at: <https://doi.org/10.1101/321729>.
43. Röhling, S., Linne, A., Schellhorn, J., Hosseini, M., Dencker, T., & Morgenstern, B. (2020). The number of k-mer matches between two DNA sequences as a function of k and applications to estimate phylogenetic distances. *Plos one*, 15(2), e0228070.
44. Chen, Z., Erickson, D. L., & Meng, J. (2020). Benchmarking hybrid assembly approaches for genomic analyses of bacterial pathogens using Illumina and Oxford Nanopore sequencing. *BMC genomics*, 21(1), 1-21.
45. Baxevanis, A.D., 2020. Biological sequence databases. *Bioinformatics*. 4th Edition ed. New York: John Wiley & Sons, pp.1-18.
46. Kim, J., Ji, M., & Yi, G. (2020). A Review on Sequence Alignment Algorithms for Short Reads Based on Next-Generation Sequencing. *IEEE Access*, 8, 189811-189822.
47. Venuto, D. and Bourque, G., 2018. Identifying co-opted transposable elements using comparative epigenomics. *Development, growth & differentiation*, 60(1), pp.53-62.
48. Manchanda, N., Portwood, J.L., Woodhouse, M.R., Seetharam, A.S., Lawrence-Dill, C.J., Andorf, C.M. and Hufford, M.B., 2020. GenomeQC: a quality assessment tool for genome assemblies and gene structure annotations. *BMC genomics*, 21(1), pp.1-9.
49. Yin, C., 2017. Encoding DNA sequences by integer chaos game representation. *arXiv preprint arXiv:1712.04546*.

50. Xiaolei, W., Wubin, Q., Chenggang, Z. and Dongsheng, Z., 2015. Kmer-indexer: A Fast K-mer Indexing Program. *Studies in health technology and informatics*, 216, pp.1083-1083
51. Brodsky, L., Kogan, S., BenJacob, E., & Nevo, E. (2010). A binary search approach to whole-genome data analysis. *Proceedings of the National Academy of Sciences*, 107(39), 16893-16898.
52. Pingali, K.D., Tanay, K.C.P. and Baruah, P.K., 2017, March. GPU accelerated suffix array construction for large genome sequences. In *Innovations in Information, Embedded and Communication Systems (ICIIECS)*, 2017 International Conference on (pp. 1-6). IEEE.
53. Bayat, A., Deshpande, N.P., Wilkins, M.R. and Parameswaran, S., 2018. Fast Short Read De-Novo Assembly Using Overlap-Layout-Consensus Approach. *IEEE/ACM transactions on computational biology and bioinformatics*, 17(1), pp.334-338.
54. Ekblom, R. and Wolf, J.B., 2014. A field guide to whole-genome sequencing, assembly and annotation. *Evolutionary applications*, 7(9), pp.1026-1042.
55. Kulkarni, P. and Frommolt, P., 2017. Challenges in the Setup of Large-scale Next-Generation Sequencing Analysis Workflows. *Computational and Structural Biotechnology Journal*, 15, pp.471-477.
56. Pan, T., Flick, P., Jain, C., Liu, Y. and Aluru, S., 2017. Kmerind: A flexible parallel library for k-mer indexing of biological sequences on distributed memory systems. *IEEE/ACM transactions on computational biology and bioinformatics*.
57. Shafin, K., Pesout, T., Lorig-Roach, R., Haukness, M., Olsen, H. E., Bosworth, C., ... & Paten, B. (2020). Nanopore sequencing and the Shasta toolkit enable efficient de novo assembly of eleven human genomes. *Nature biotechnology*, 38(9), 1044-1053.
58. Giordano, F., Stammnitz, M. R., Murchison, E. P., & Ning, Z. (2018). scanPAV: a pipeline for extracting presence-absence variations in genome pairs. *Bioinformatics*, 34(17), 3022-3024.
59. Mikheenko, A., Prjibelski, A., Saveliev, V., Antipov, D., & Gurevich, A. (2018). Versatile genome assembly evaluation with QUAST-LG. *Bioinformatics*, 34(13), i142-i150.
60. Rhie, A., Walenz, B. P., Koren, S., & Phillippy, A. M. (2020). Merqury: reference-free quality, completeness, and phasing assessment for genome assemblies. *Genome biology*, 21(1), 1-27.
61. Chen, Y., Zhang, Y., Wang, A. Y., Gao, M., & Chong, Z. (2021). Accurate long-read de novo assembly evaluation with Inspector. *Genome biology*, 22(1), 1-21.
62. Castro, C. J., & Ng, T. F. F. (2017). U50: a new metric for measuring assembly output based on non-overlapping, target-specific contigs. *Journal of Computational Biology*, 24(11), 1071-1080.
63. Taylor, A.S., Barros, D., Gobet, N., Schuepbach, T., McAllister, B., Aeschbach, L., Randall, E.L., Trofimenko, E., Heuchan, E.R., Barszcz, P. and Ciosi, M., 2022. Repeat Detector: versatile sizing of expanded tandem repeats and identification of interrupted alleles from targeted DNA sequencing. *NAR Genomics and Bioinformatics*, 4(4), .lqac089. Available at: <https://doi.org/10.1093/nargab/lqac089>.

- 739 64. Liao, X., Li, M., Hu, K., Wu, F. X., Gao, X., & Wang, J. (2021). A sensitive repeat
740 identification framework based on short and long reads. *Nucleic acids research*, 49(17), e100-
741 e100.
- 742 65. Guiglielmoni, N., Houtain, A., Derzelle, A., Van Doninck, K., & Flot, J. F. (2021).
743 Overcoming uncollapsed haplotypes in long-read assemblies of non-model organisms. *BMC*
744 *bioinformatics*, 22(1), 1-2

Box 1(on next page)

Box 1 N50 Metric Calculation

```

1  Let C array of scaffolds lengths C={c1, c2, c3,...,ci}
2  Let AL is Sum C[]
3  Let AL50= Sum C[]/2
4  Let incL is the incremental length with initial value incL=0
5  Let CN50[] is an empty array to store N50 scaffolds' members
6  Let i is the length of array C[]
7  Sort C descending
8  For n=1 to n=i
9  Do
10 IF (incL= AL50)
11 Break For Loop
12 ELSE
13
14 incL= incL + C[n].Length
15 Add C[n] to CN50[]
16 End IF
17 End For
18 return CN50[]
19 Thus:
20 N50=min( CN50[])
21
22
23

```

Box 1 N50 Metric Calculation

Box 2(on next page)

Index hit represent the overlapping computational time complexity

Let $N=5$ reads
Let All-against-all $O(N)^2 = 25$ times
Current index hit = 10
Thus: Overlapping computational time complexity is $O(2N)$ which is
liner algorithm.

Box 2 Index hit represent the overlapping computational time complexity

Box 3(on next page)

Box 3 Search string

(Genome assembly OR genome analysis OR indexing genome data OR genome misassembly OR parallelize genome data OR repetitive genome sequence identification) AND (accuracy OR performance OR enhancement) AND (method OR approach) AND (De novo approach OR reference guided approach OR hybrid approach) AND language (English)

Box 3 Search string

Figure 1

Figure 1 Biological activities of TE and TR change genome size

Transposon elements copy themselves
into another genome region



Tandem repeats increase their copy number

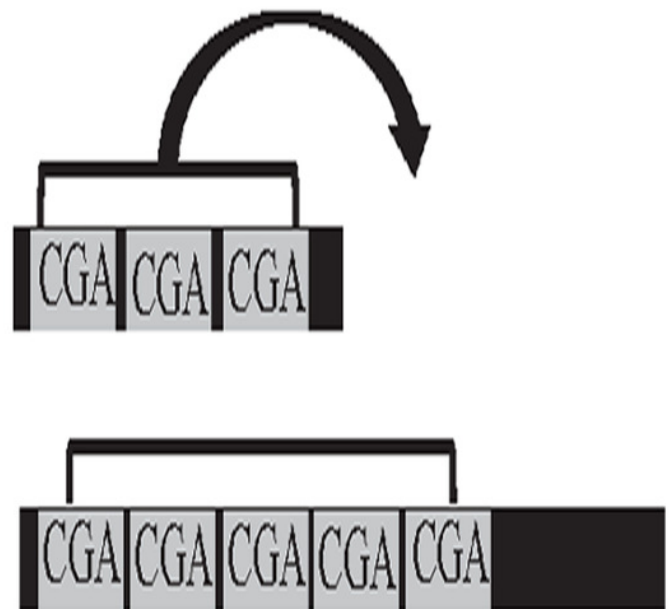


Figure 2

Figure 2 String exact and approximate match

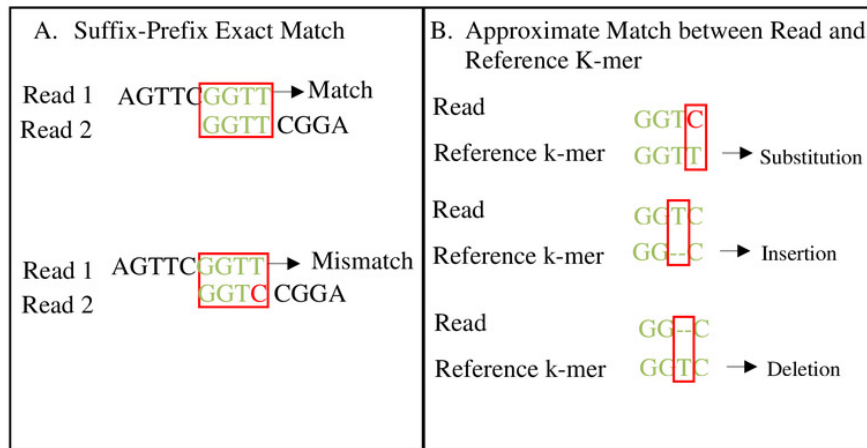


Figure 2

String exact and approximate match

Figure 3

Figure 3 Overlapping suffix-prefix match

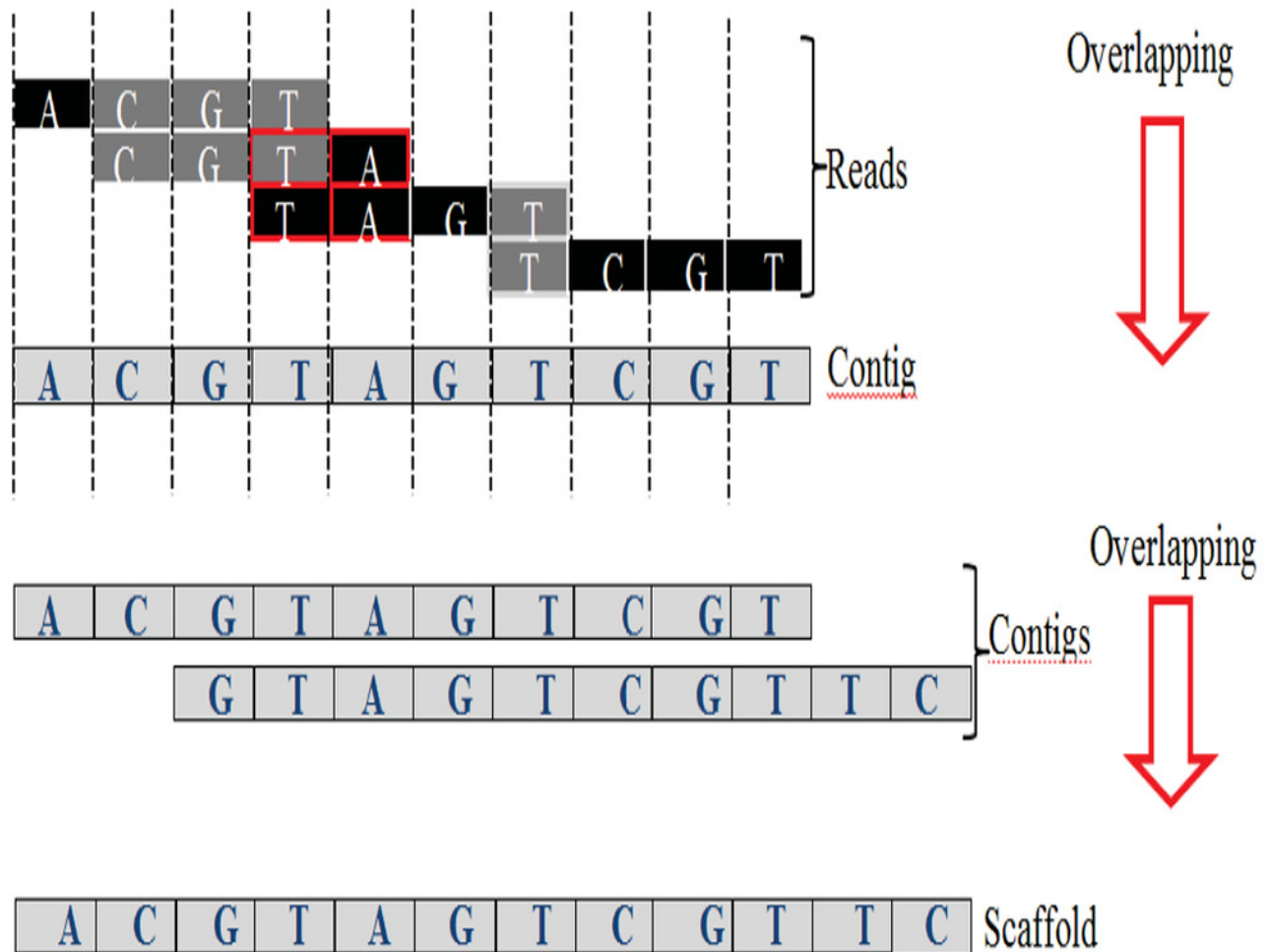


Figure 4

Figure 4 The impact of overlapping length on assembly accuracy

Given Genome (G): ACTTCA**CGTCGTCGTCG**TTGATCAA

Long Overlapping Length

Read (r1): ACTTCA**CGTCGTCG** ← Overlapping distance=6

Read r2: **TCGTCG**TCGTTGATCAA



Assembly Overlapping

Genome G: ACTTCA**CGTCGTCGTCG**TTGATCAA **Correct Genome**

Short Overlapping Distance

Read r1: ACTTCA**CGTCGTCG** ← Overlapping distance=3

Read r2: **TCGTCG**TCGTTGATCAA



Assembly Overlapping

Genome G: ACTTCA**CGTCGTCGTCGTCGTCG**TTGATCAA
Wrong Genome (Mis-assembly)

Figure 5

Figure 5 Reads mapping against reference genome (reads alignment)

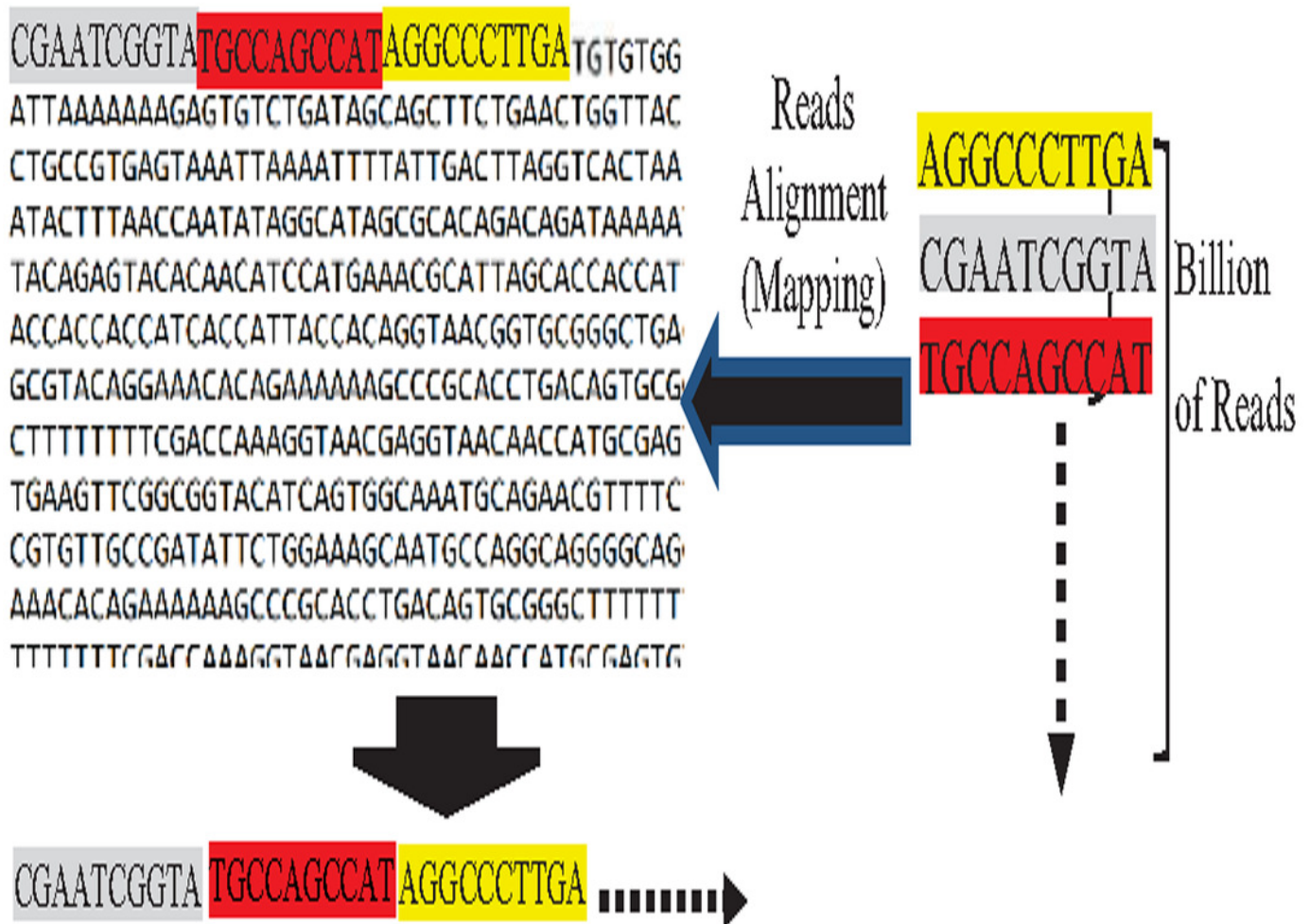


Figure 6

Figure 6 Sequencing coverage depth

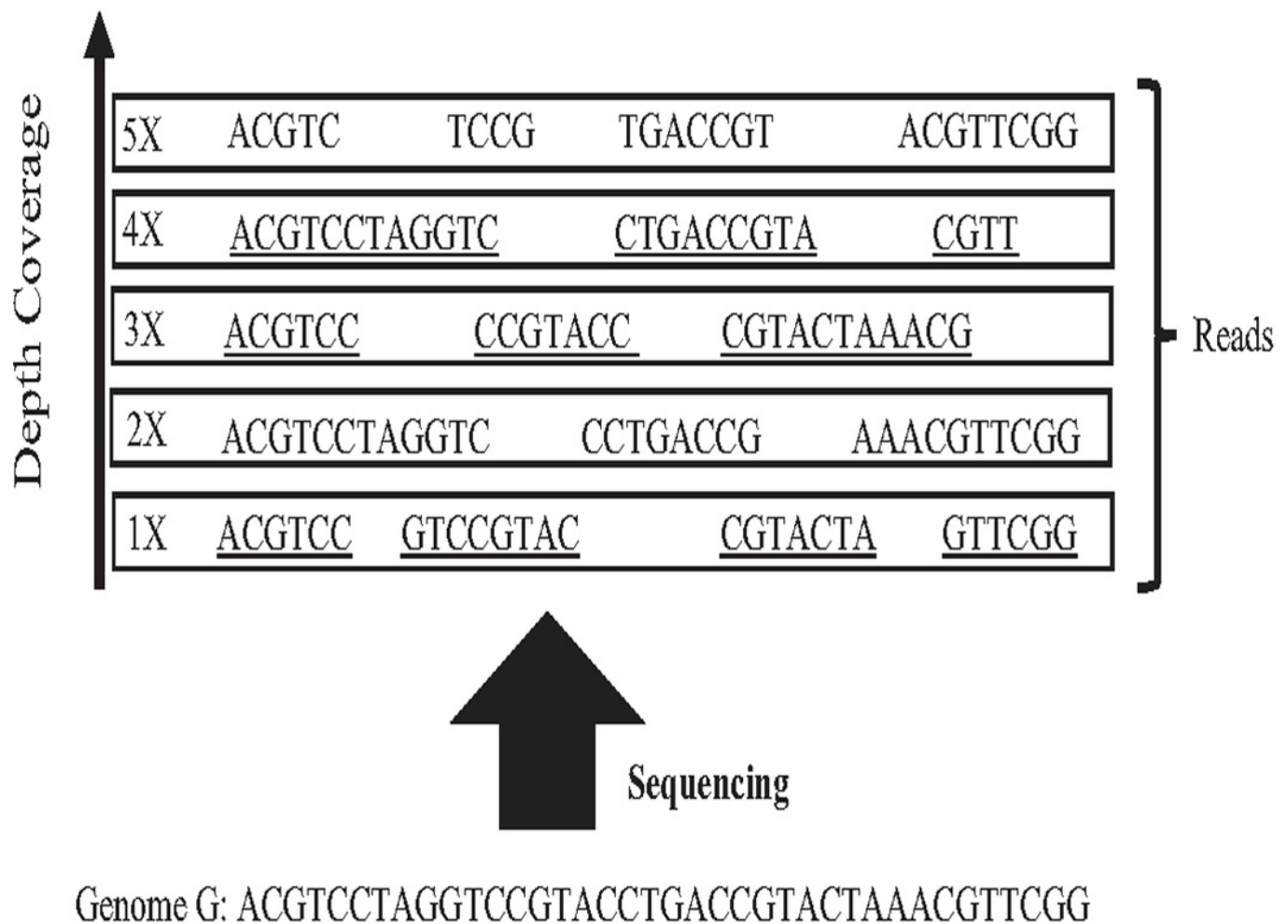


Figure 7

Figure 7 Mis-assembly in de novo approach

Original Genome (G): TGGGACTGG

Sequencing

Input from biology to computer Science:

$S = \{ \text{GACT}, \text{ACTGG}, \text{TGGGAC} \}$

Overlapping algorithm

Output from computer science to biology:

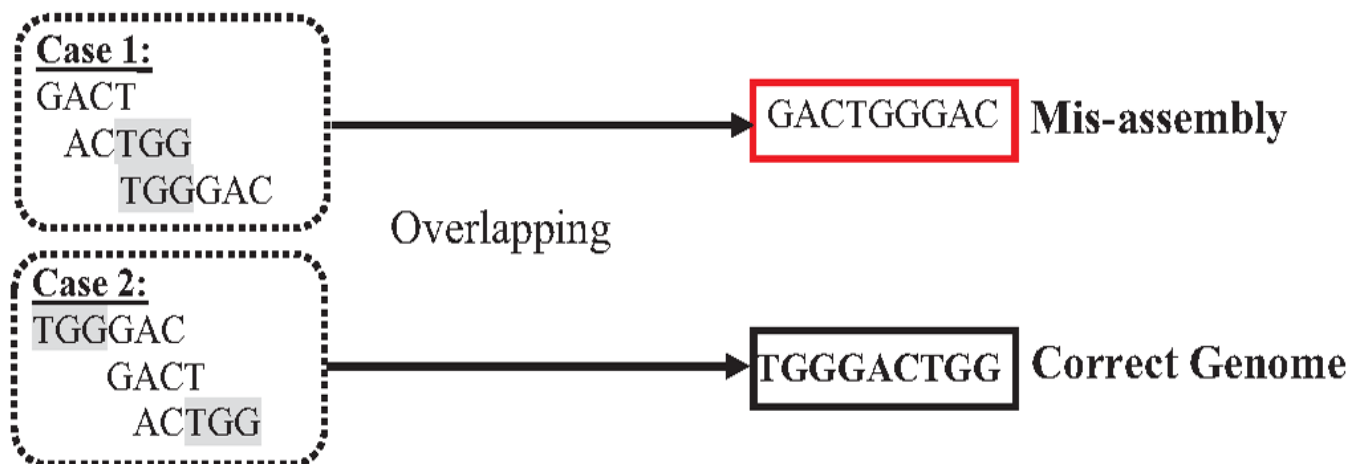


Figure 8

Figure 8 Mis-assembly in greedy algorithm

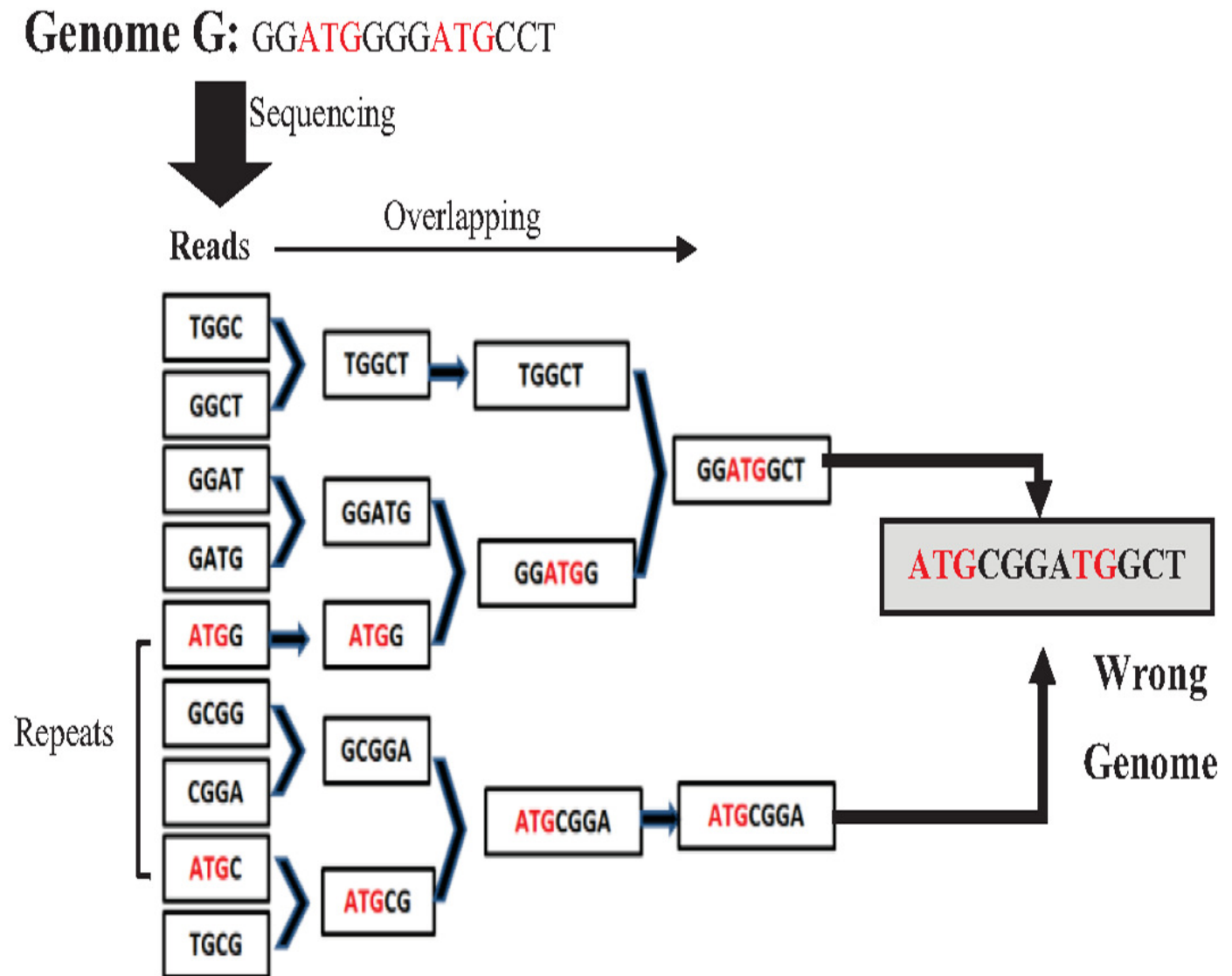
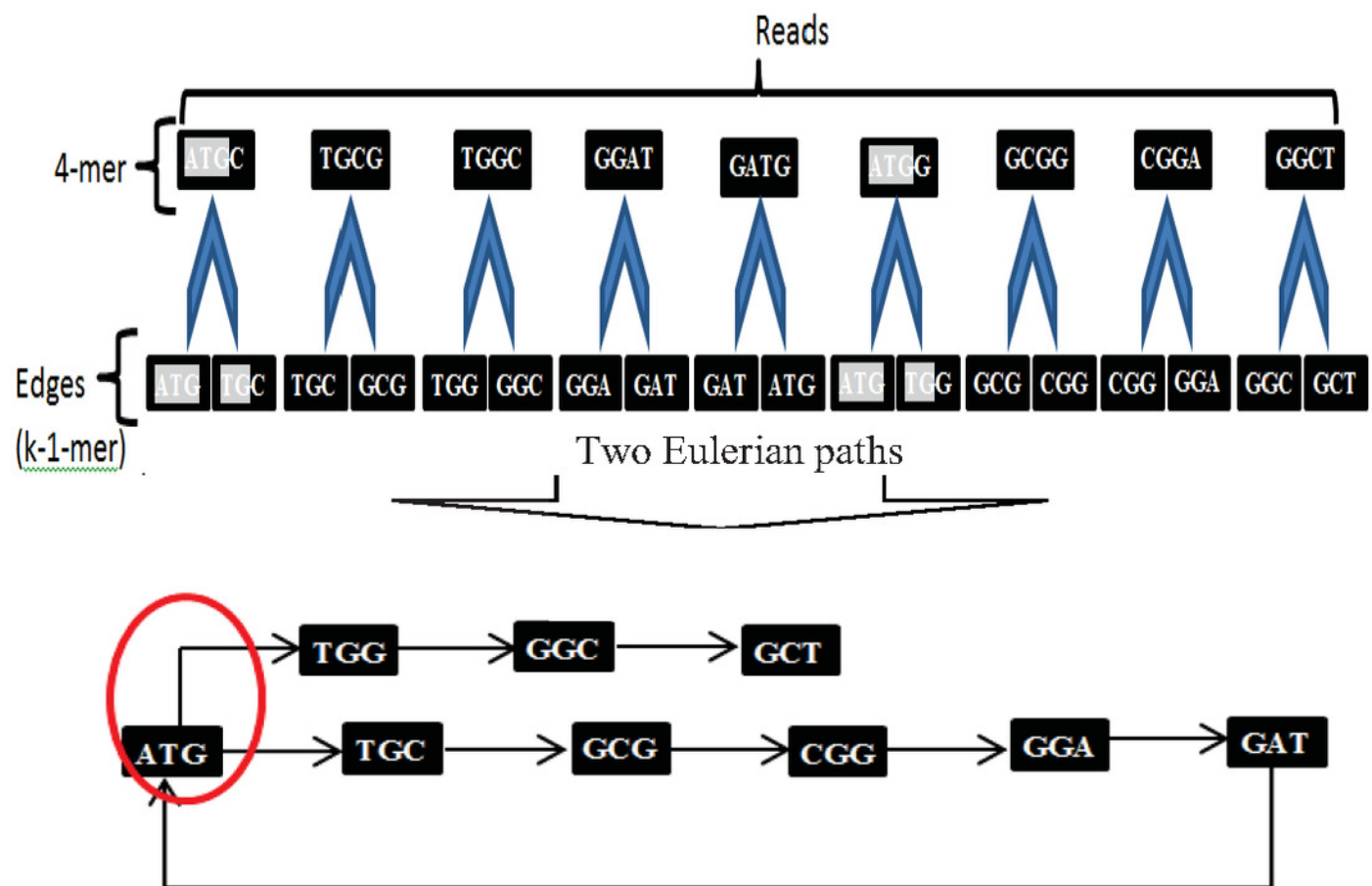


Figure 9

Figure 9 Fragmented assembly in DBG algorithm



First assembly: ATGGCT

Fragmented Assembly Second assembly: ATGCGGAT

Figure 10

Figure 10 Mis-assembly in reference-guided approach

Input from biology to computer Science:

$R = \{GGATGCGATGGCT\}$ Reference Genome-----Input 1

$R = \{ATGCG, GAGTT, ATGCA\}$ read to be aligned-----Input 2

Input String: **ATGCG, GAGTT, ATGCA**

Reads alignment

Reference Genome (G): **ATGCGAGGCAGAGTT**

Case 1 **ATGCGATGCAGAGTT** *Correct Genome*

Case 2 **ATGCAATGCGGAGTT** *Wrong Genome*
Mis-Assembly

Output from computer science to biology (if case 2):

ATGCAATGCGGAGTT which is wrong genome

Figure 11

Figure 11 Overlapping and reads alignment time complexity

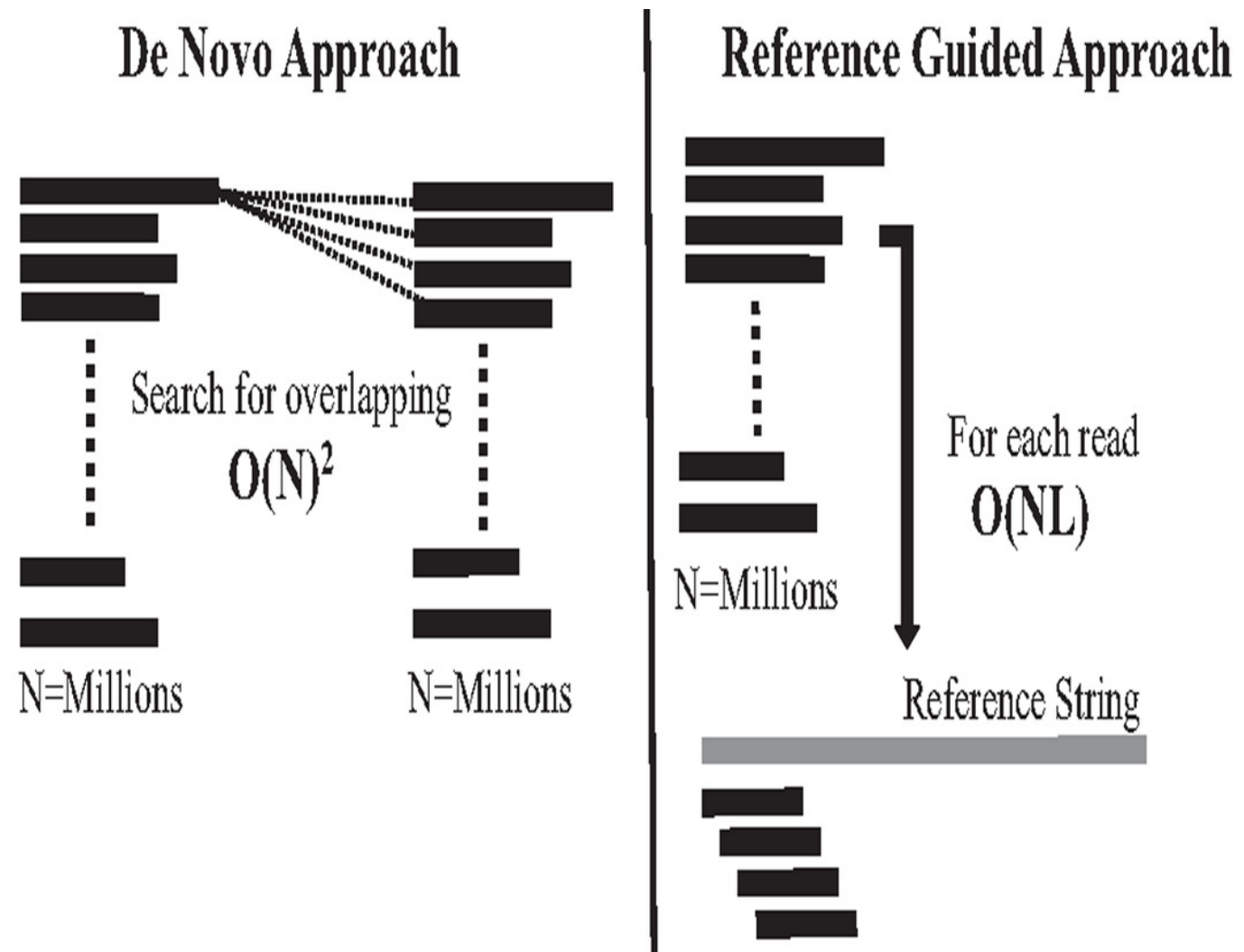


Figure 12

Figure 12 Search Map

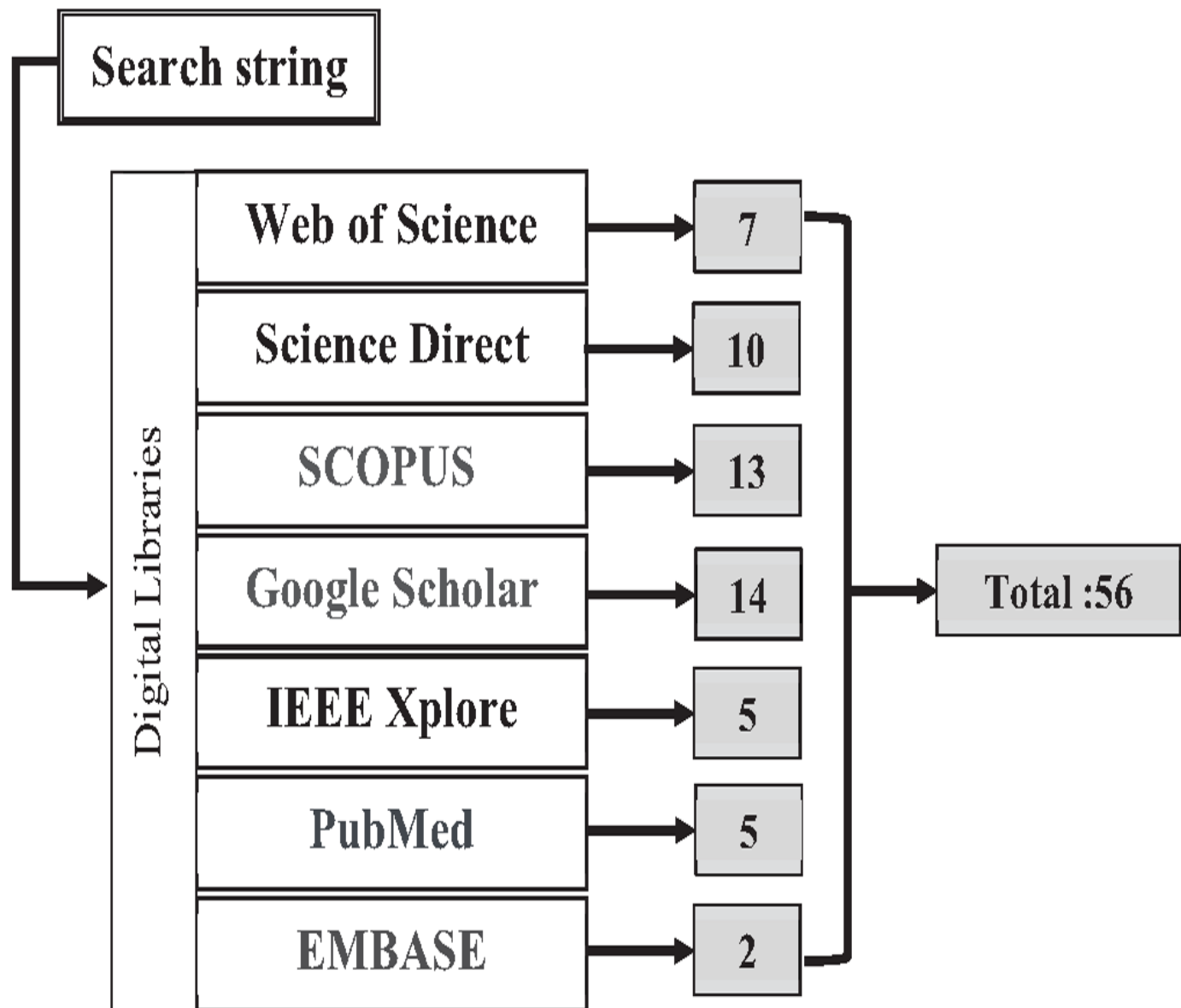


Figure 13

Figure 13 Repetitive sequences frequency greater than coverage depth

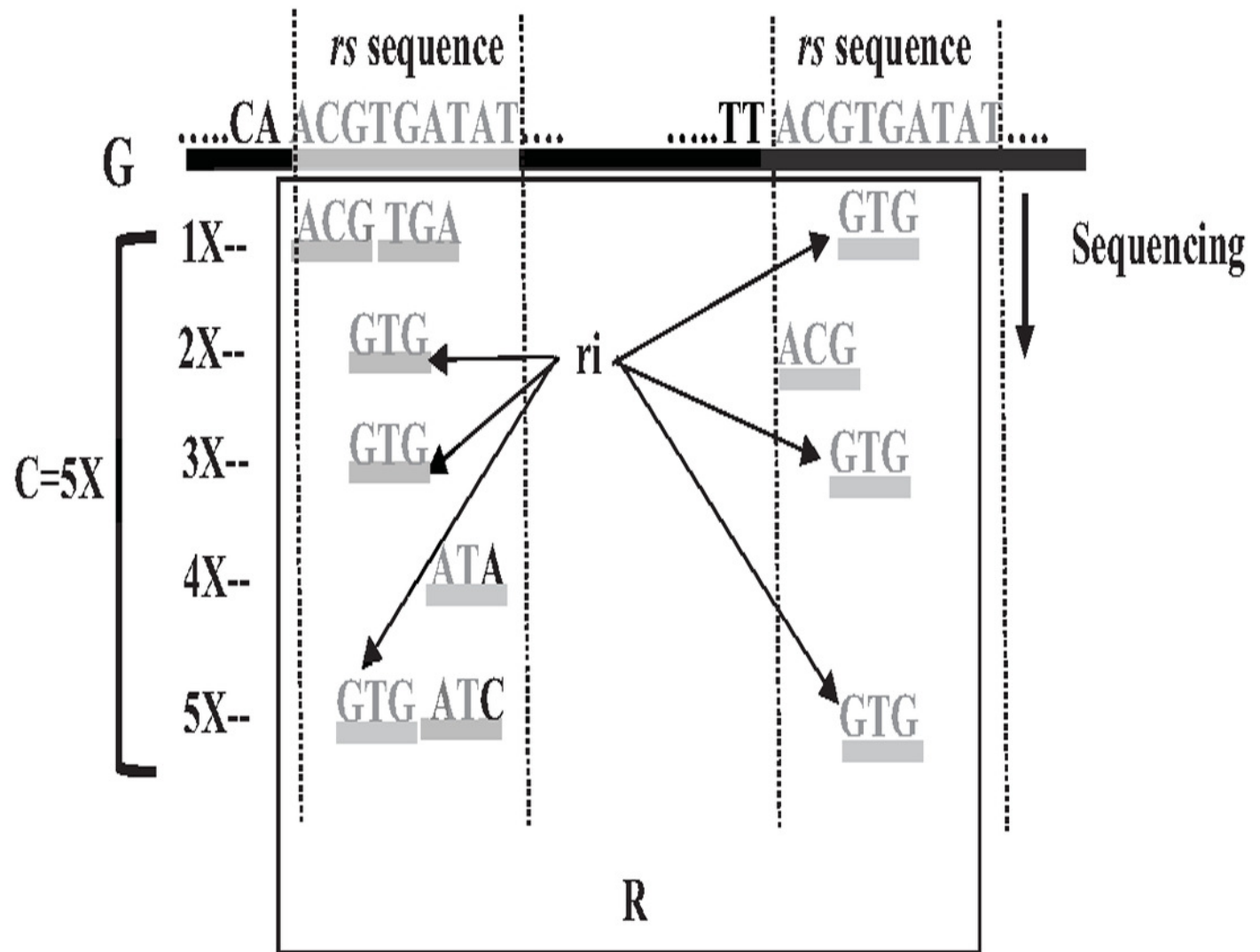


Figure 14

Figure 14 Repetitive sequence as substring of read

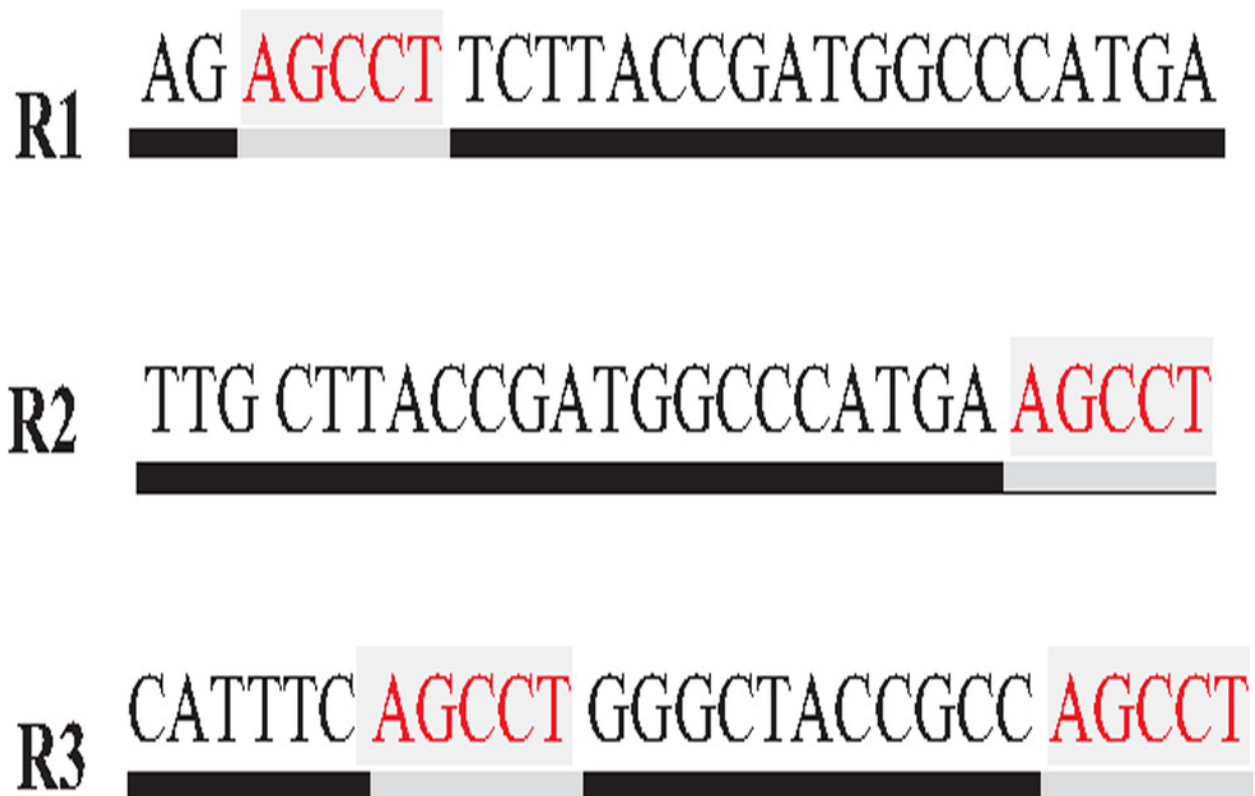


Figure 15

Figure 15 Reads' K-mers Frequency in REPdenovo

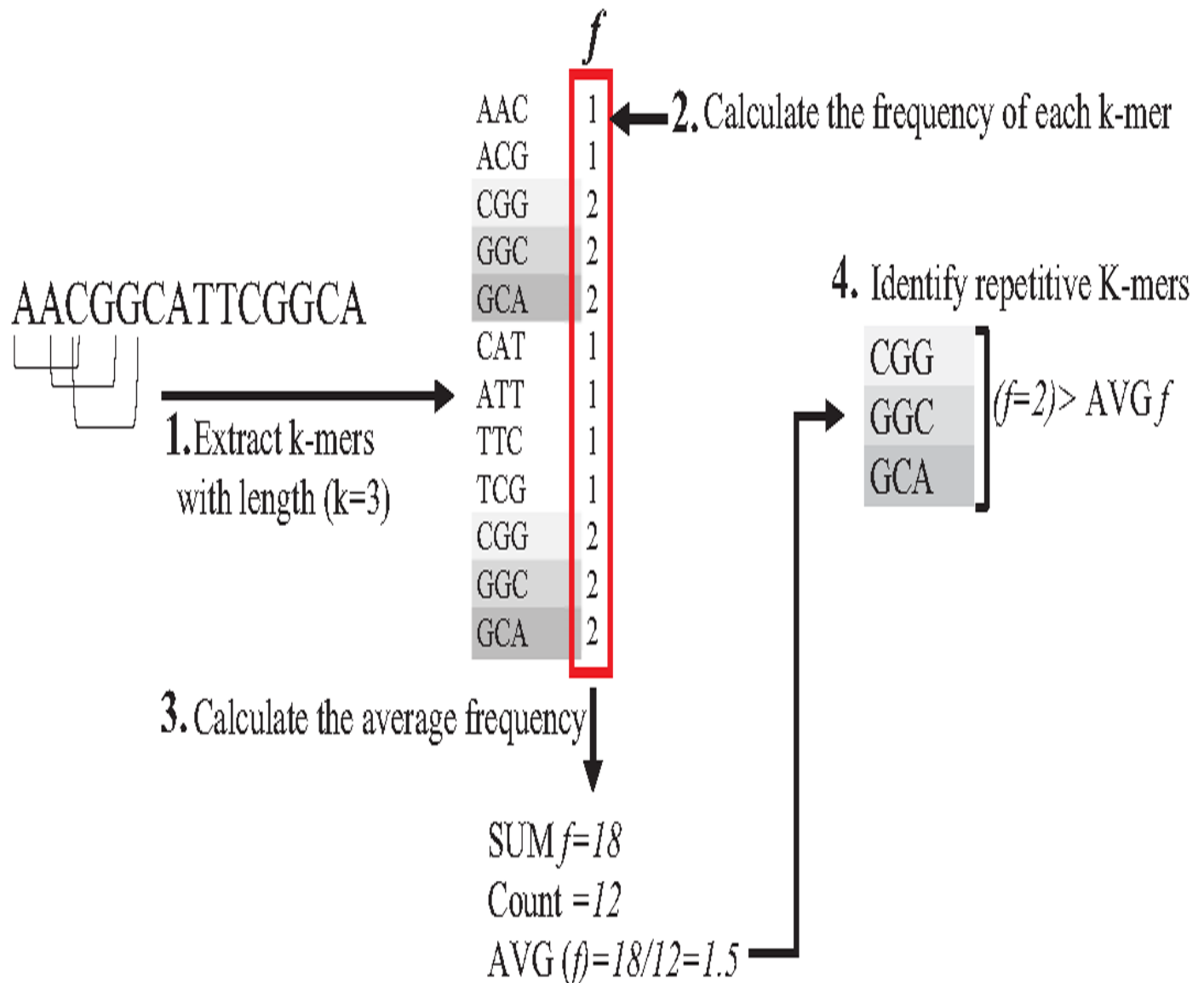
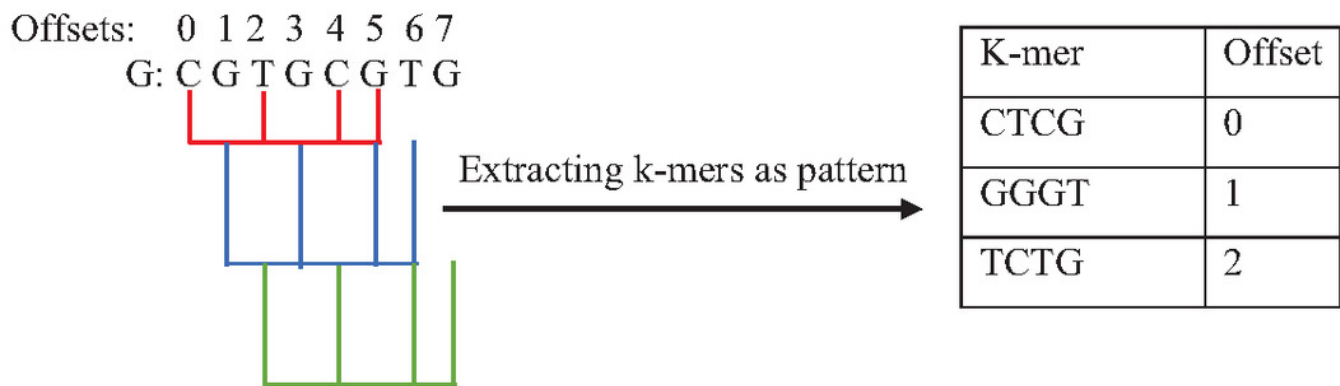
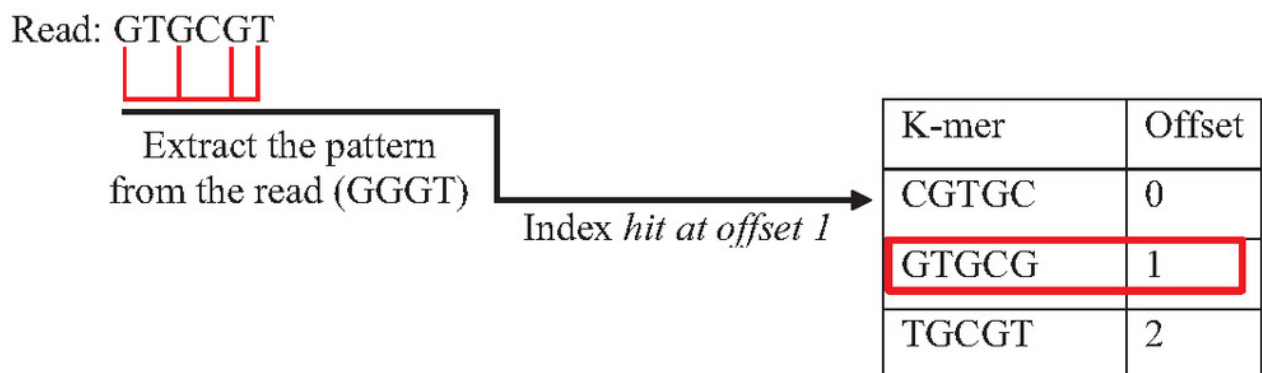


Figure 16

Figure 16 Querying k-mer index with pattern match



Query the previous index:



Verification:

Offset: 0 1

G: CGTGC GTG

R: GTGCGT

Figure 17

Figure 17 Binary search to querying k-mer index

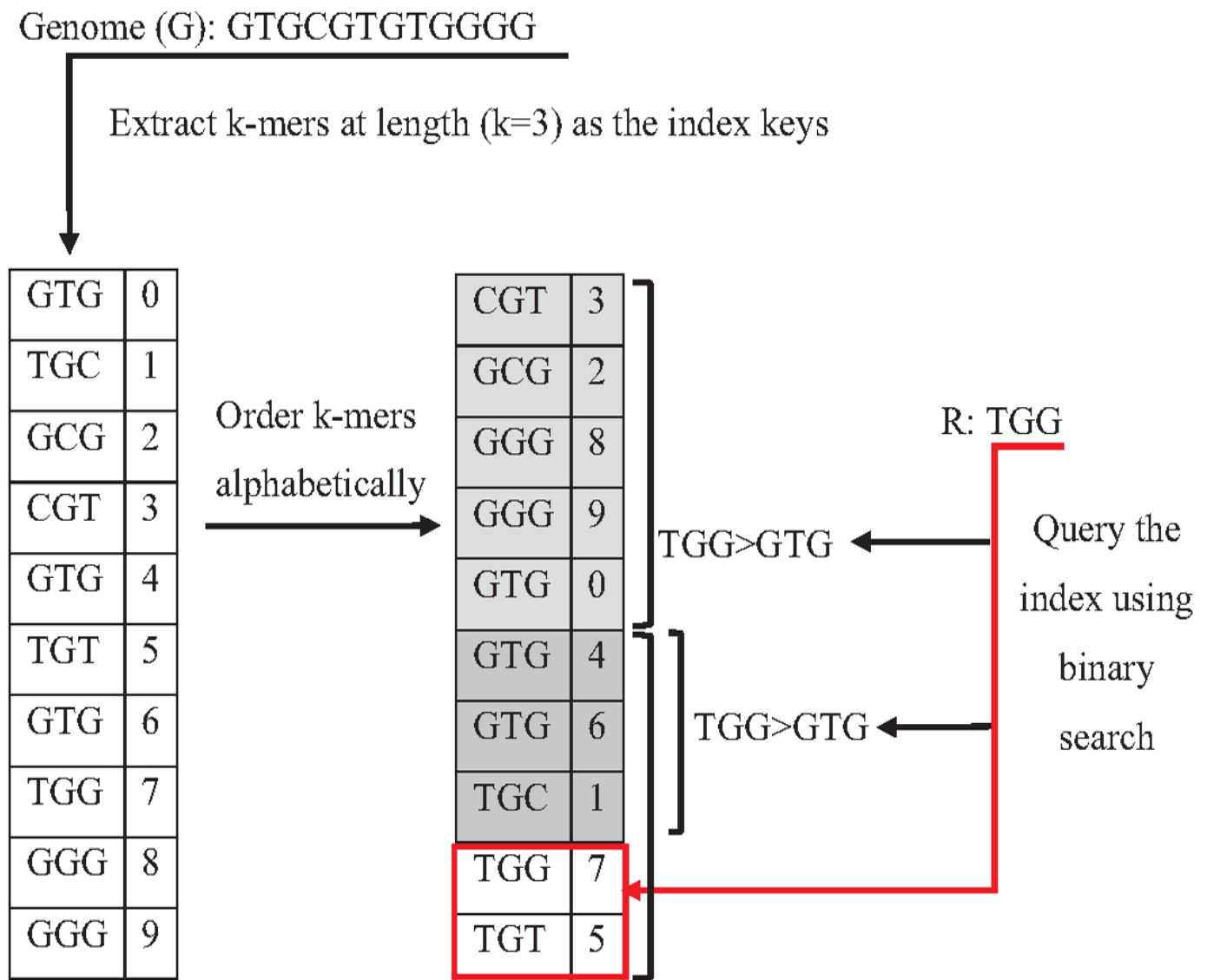


Figure 18

Figure 18 Construction of hash index and collision resulted by repeat

Genome (G): GTGCGTGTGGGGG

Index the k-mers with length k=3 in to groups

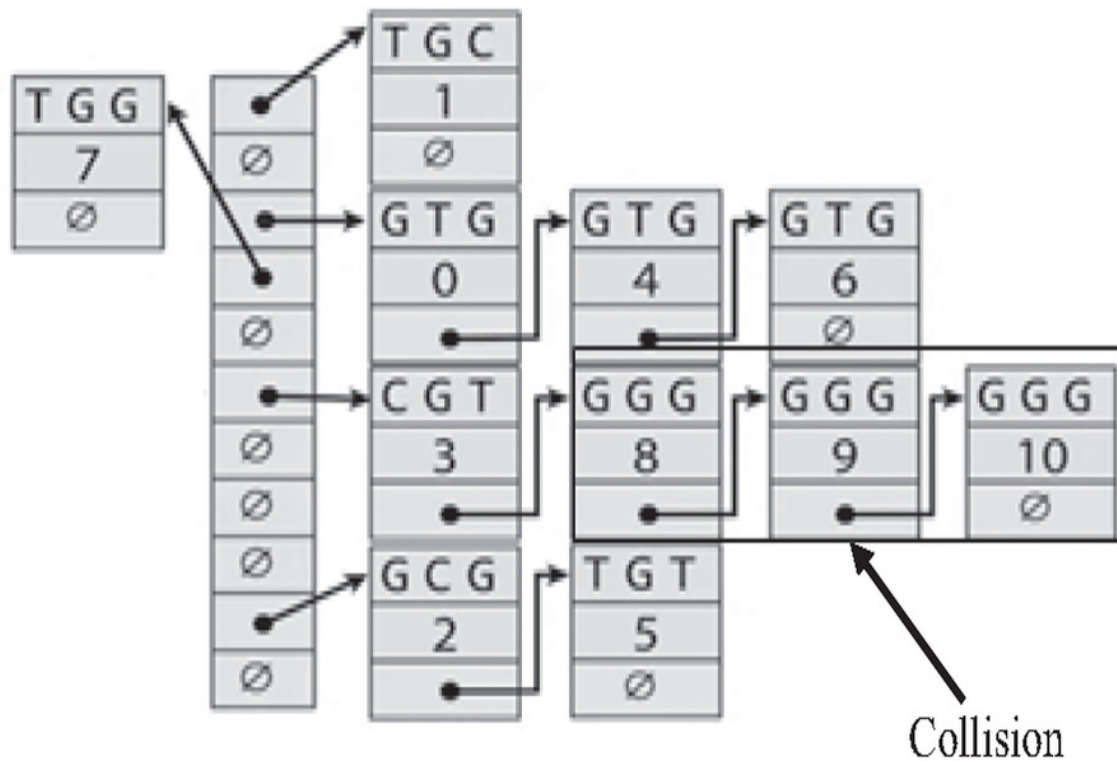


Figure 19

Figure 19 Prefix-tree index construction

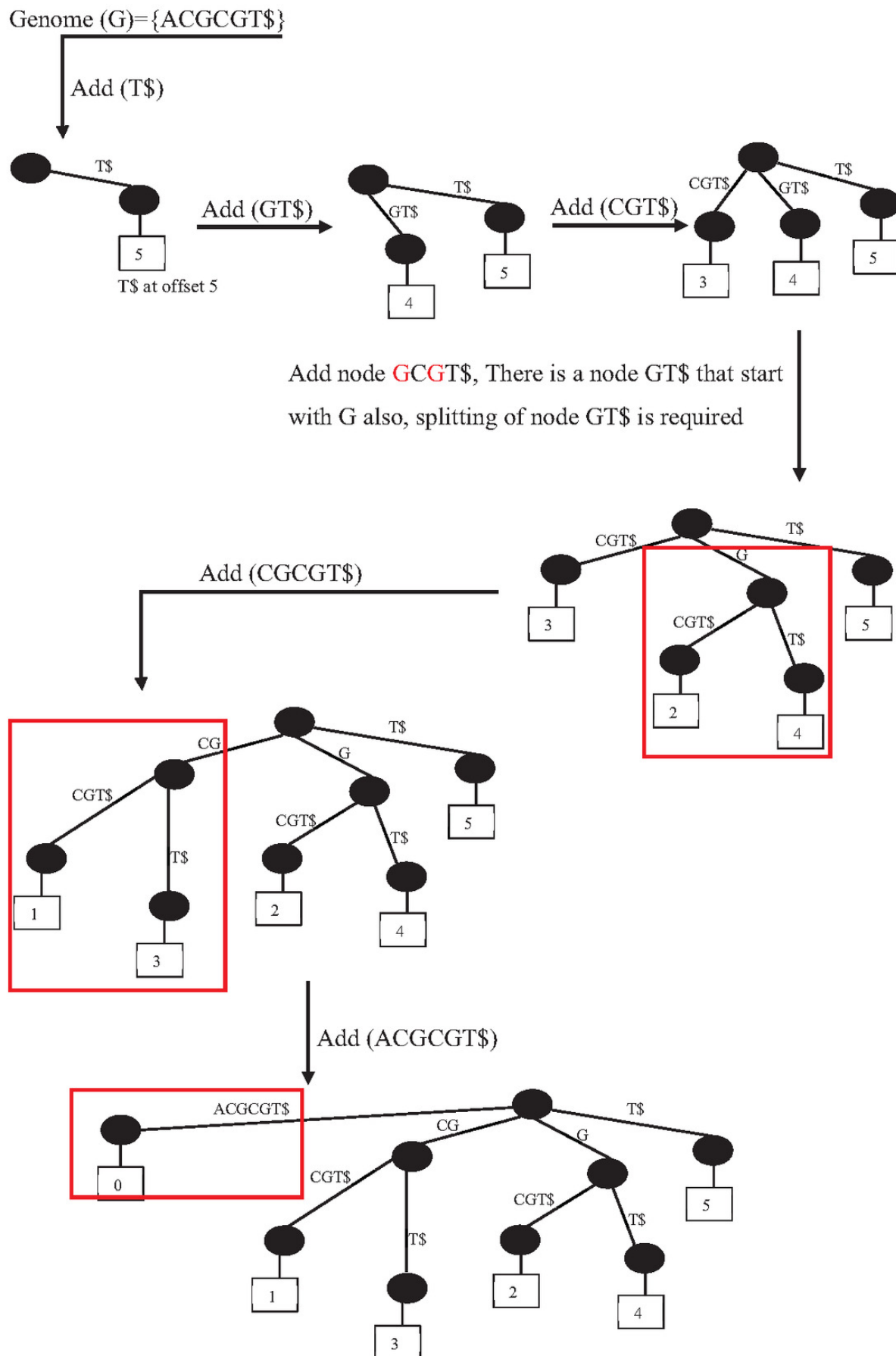


Figure 20

Figure 20 Pigeonhole principle in overlapping

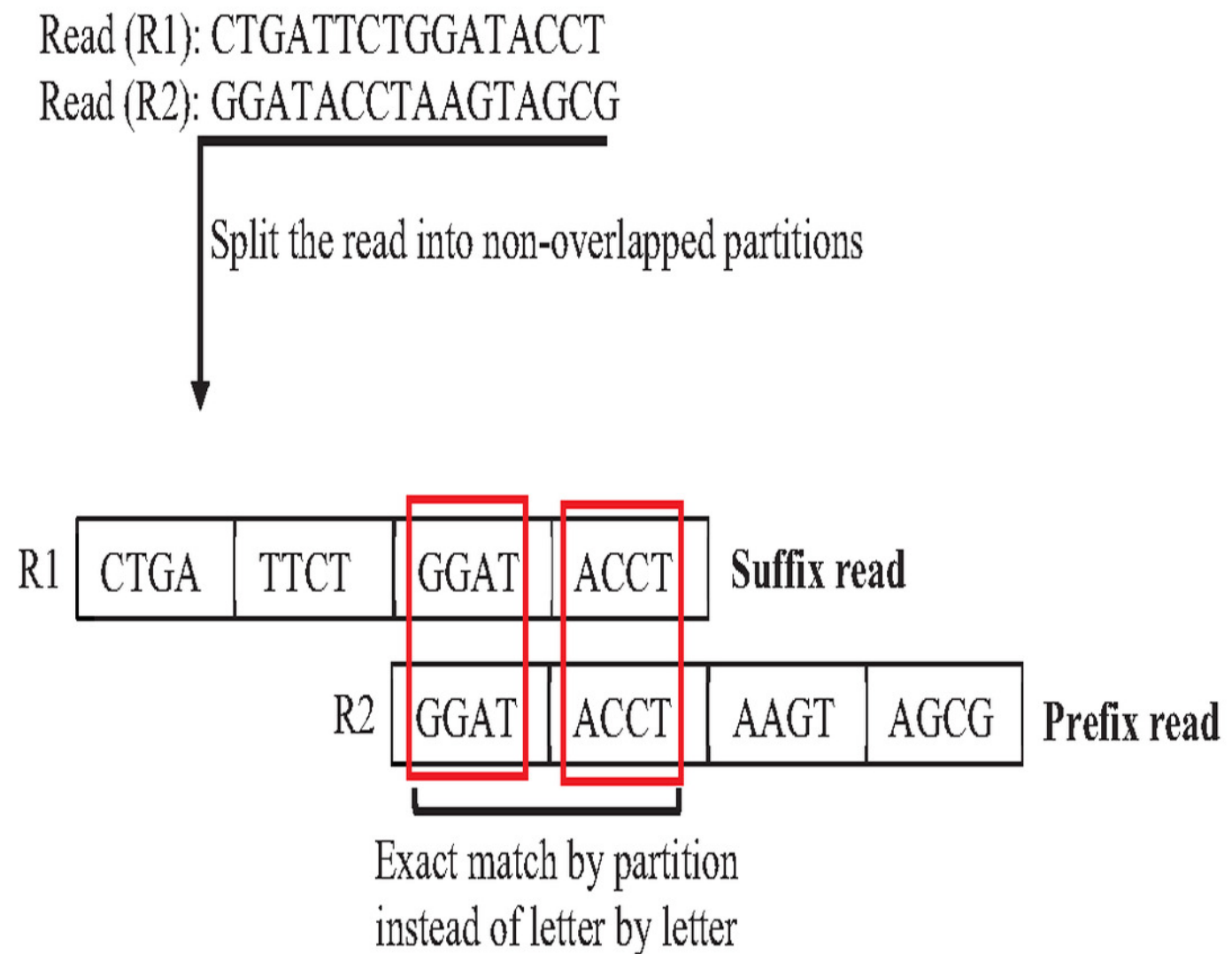


Figure 21

Figure 21 Solution to propose accurate and optimized assembly hybrid approach

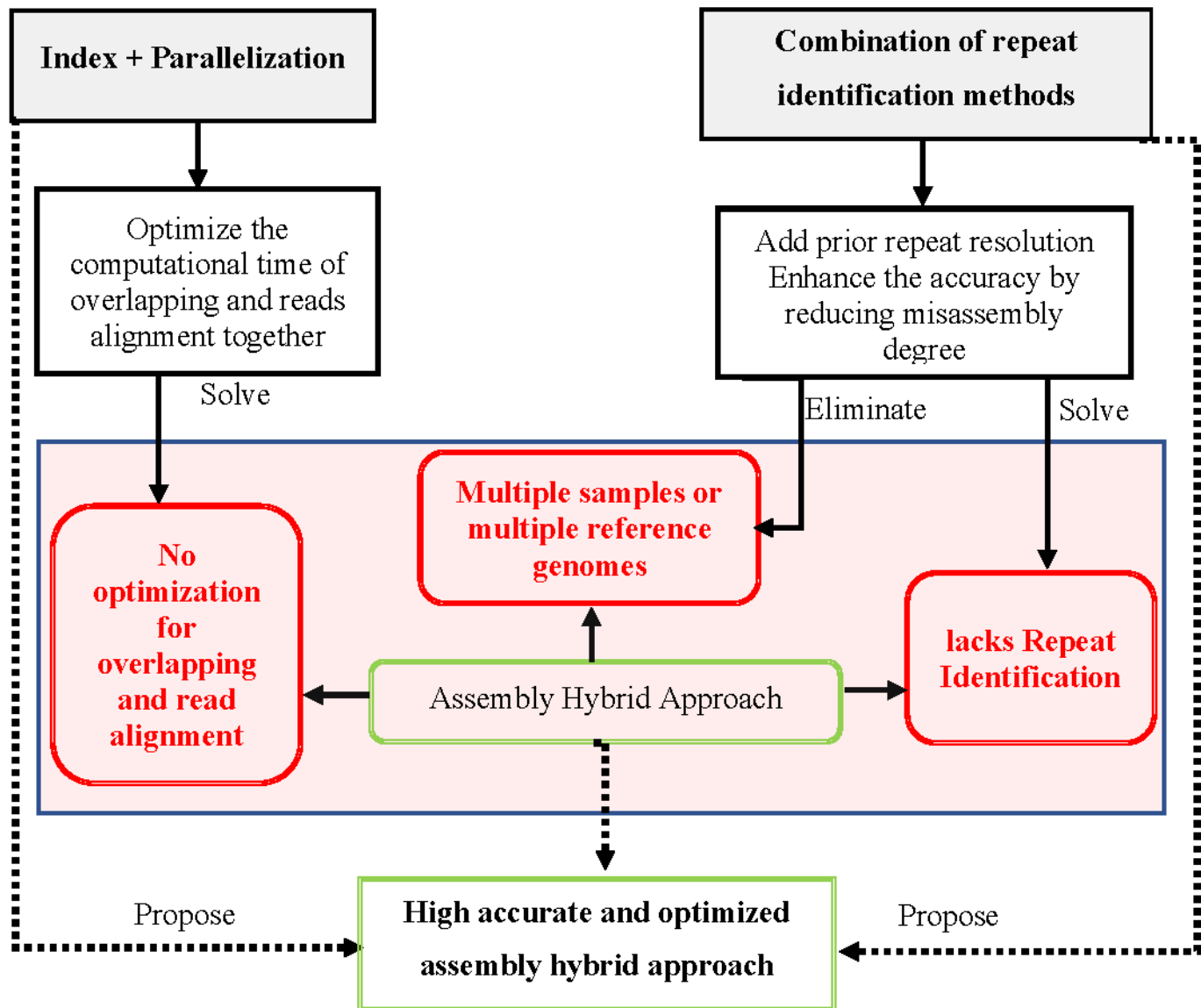


Table 1 (on next page)

Table 1 The limitations of current repeat identification methods

Table 1 The limitations of current repeat identification methods

Repeat Identification Method	Limitations
SWA	Cannot detect repetitive sequence shorter than read's length
Detection in DACCOR assembly	Identify repetitive sequences from the reference genome that are slightly different from the target genome resulting in low identification accuracy.
REPdenovo	Biological info represented in read structure are lost when the reads are chopped up into k-mers
Digital signal processing	Difficult to capture a copy number of repeats and cannot capture repeats in short DNA sequences which means that many types of tandem repeats cannot be identified
Machine learning	Not useful for analysing repeat content such as transposon element for biological evolution and big performance challenge on massive genome data

Table 2 (on next page)

Table 2 Repeat identification in *Treponema pallidum* Genome using Detection method in DACCOR

Table 2 Repeat identification in *Treponema pallidum* Genome using Detection

method in

DACCOR

Detected repetitive Sequences	<i>Treponema pallidum</i>
True positives	22382
True negatives	1116478
False positives	0
False negatives	773
Accuracy(%)	99.93

Table 3(on next page)

Table 3 DACCOR Enhance the genome resolution compared to SPAdes and EAGER method

Table 3 DACCOR Enhance the genome resolution compared to SPAdes and EAGER

method

Sample	Method	<i>Contigs from repetitive regions mapped against reference genome</i>
ARI (Coverage 157X)	SPAdes (De novo)	69.5%
	EAGER (Reference-guided approach)	42.1%
	DACCOR (Hybrid approach with prior repeat Identification)	96.8%

Table 4(on next page)

Table 4 Comparison between PT, PP index in time and space

Table 4 Comparison between PT, PP index in time and space

Metric	PT	PP
Space consumption	230 MB	298 MB
Time consumption	757 second	49 second