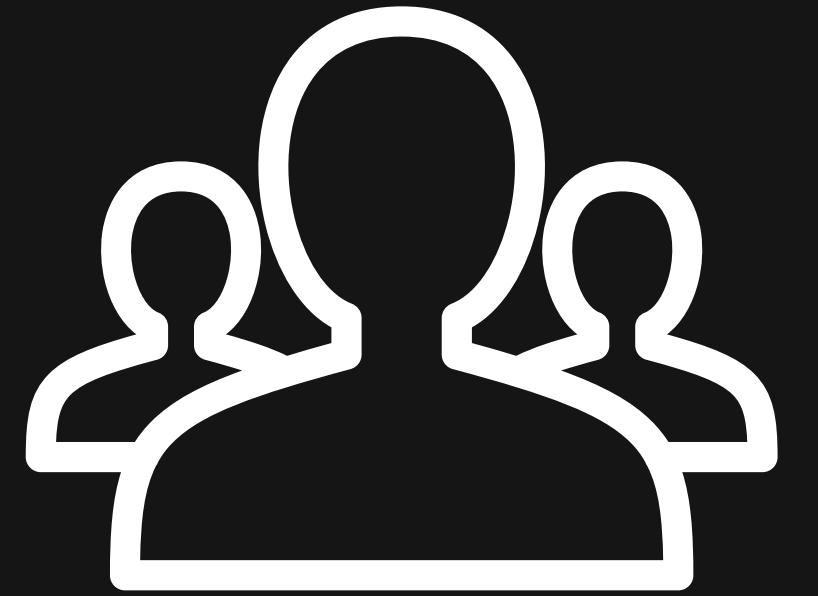




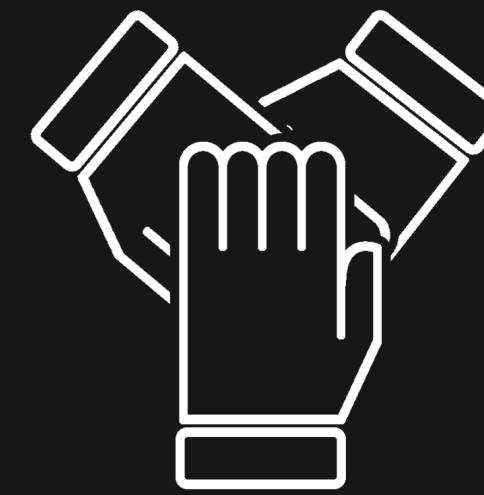
Genome Polymorphism Detection Through Relaxed de Bruijn Graph Construction



THE TEAM

GROUP - 13

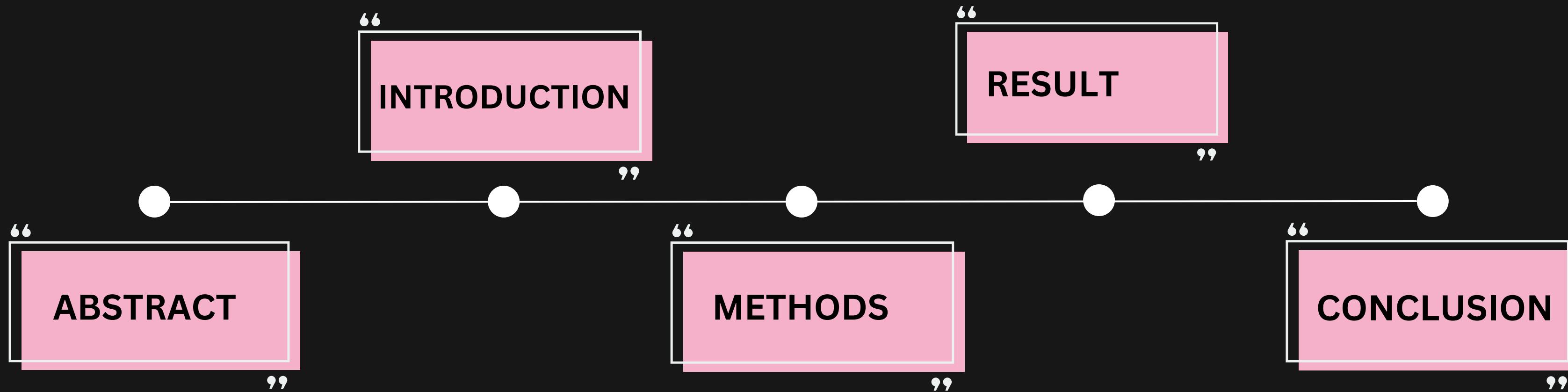
G.Venkata Vignesh
AM.EN.U4AIE21074



C.Guru Aswini Dath
AM.EN.U4AIE21023

V.Praneeth
AM.EN.U4AIE21066

TIMELINE

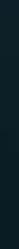




ABSTRACT

- Comparing genomes to identify polymorphisms is a difficult task, especially beyond single nucleotide polymorphisms.
- Polymorphism detection is important in disease association studies as well as in phylogenetic tree reconstruction.
- We present a method for identifying polymorphisms in genomes by using a modified version de bruijin graphs.
- Now a days, data structures are widely used in genome assembly from Next Generation Sequencing.

INTRODUCTION

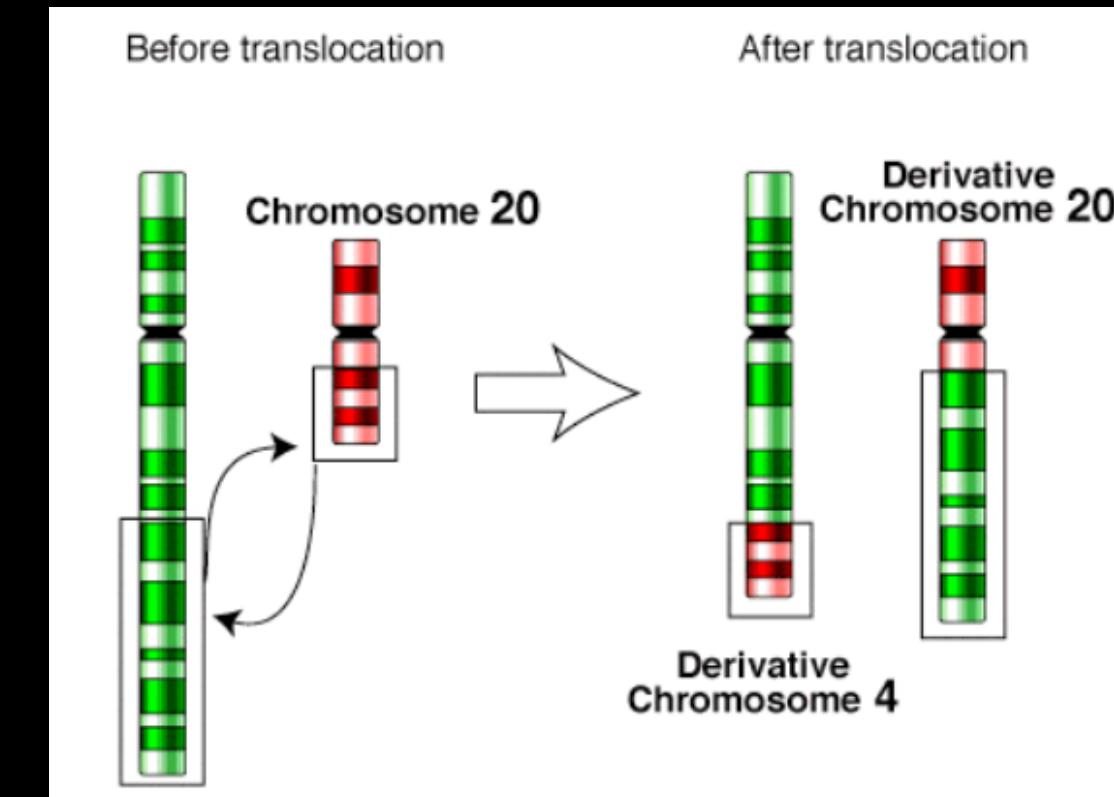




- Genomic variation such as translocations and inversions have been shown to cause many human diseases.
- Translocations have been shown to be the cause of several different types of cancer, such as Burkitt's lymphoma and acute promyelocytic leukemia.

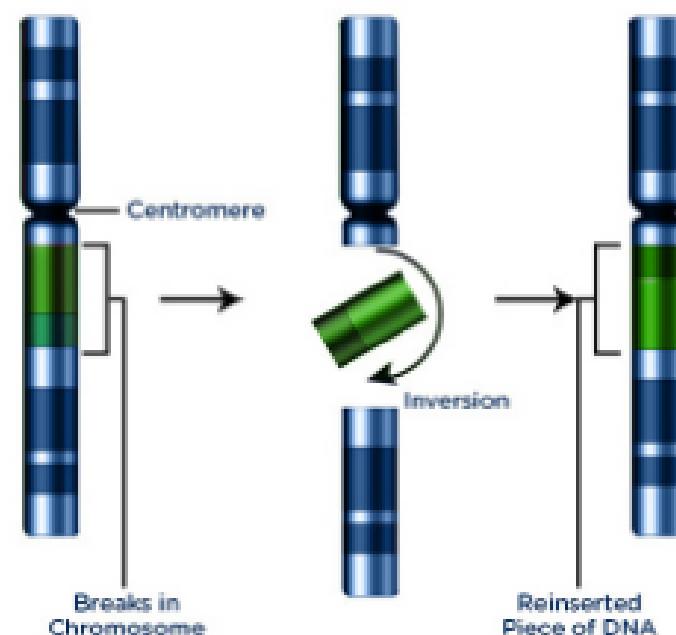
TRANSLOCATION:

A genetic change in which a piece of one chromosome breaks off and attaches to another chromosome.



Balanced Inversion

Before Inversion After Inversion

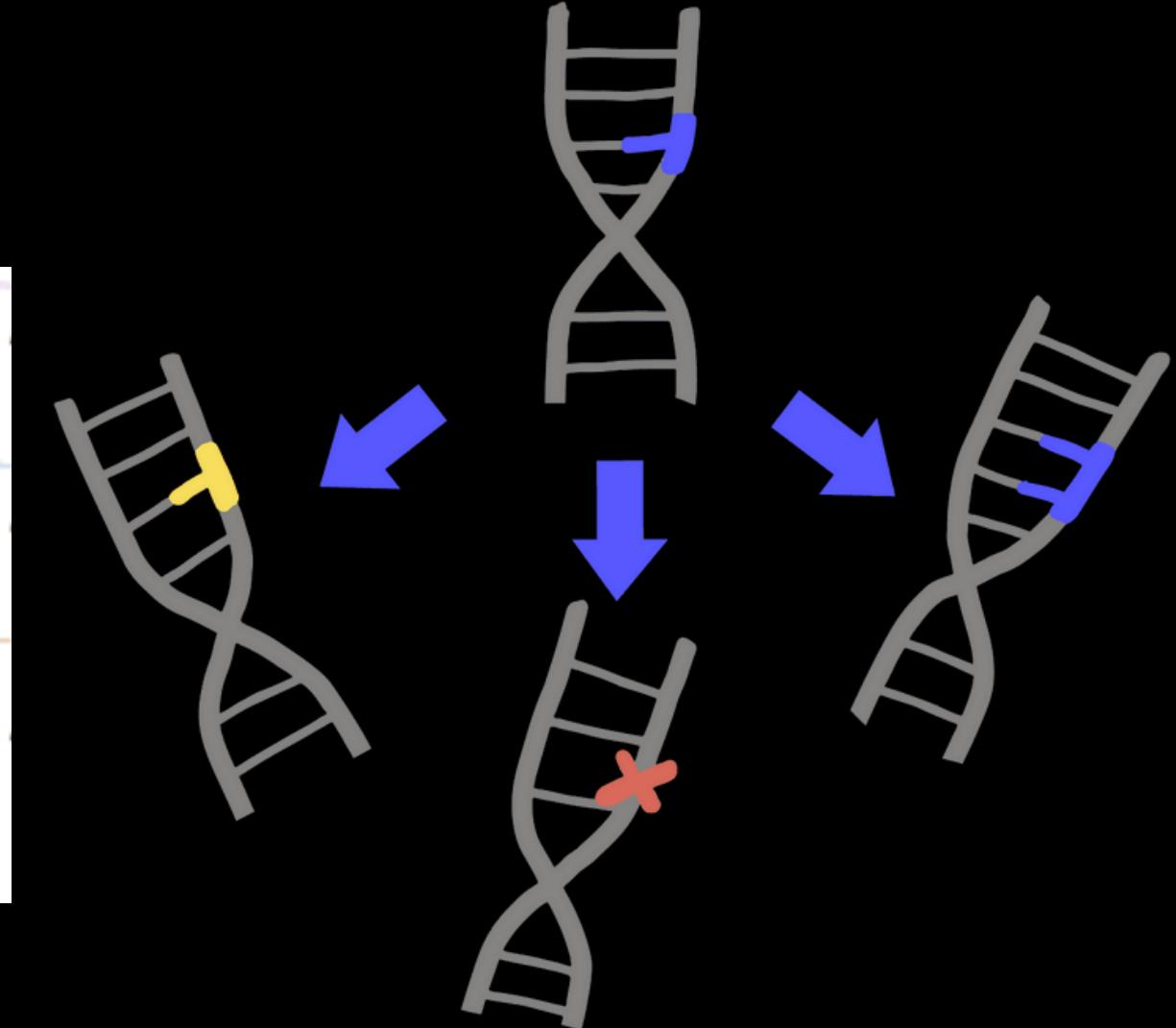


INVERSIONS:

An inversion is a chromosome rearrangement in which a segment of a chromosome becomes inverted within its original position.

Polymorphism

"Poly" many "morphe" form



POLYMORPHISM

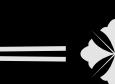
It is the presence of two or more variant forms of a specific DNA sequence that can occur among different individuals or populations.

SINGLE NUCLEOTIDE POLYMORPHISMS

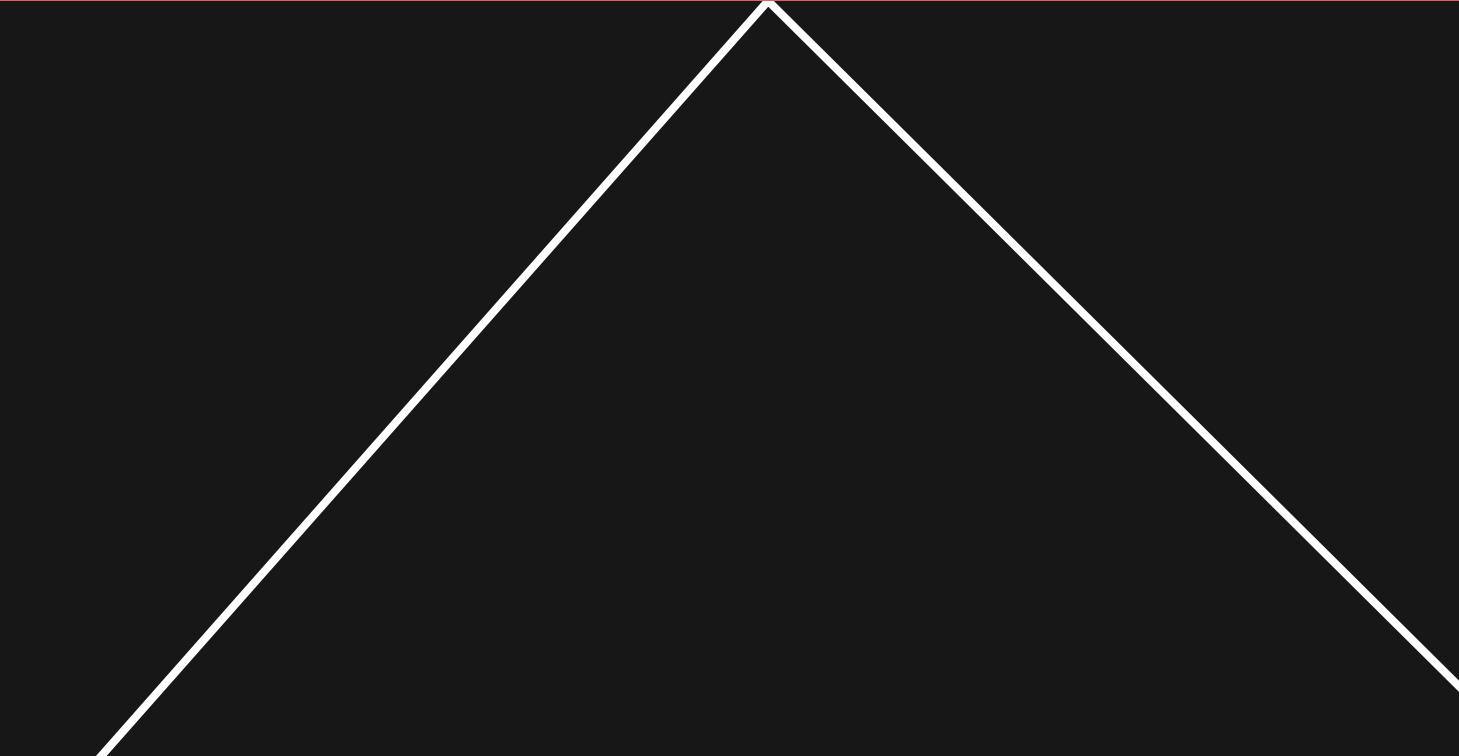
- DNA sequence variation that occurs when a single nucleotide differs between members of biological species.
- This may be substitution, insertion or deletion.

Mutation

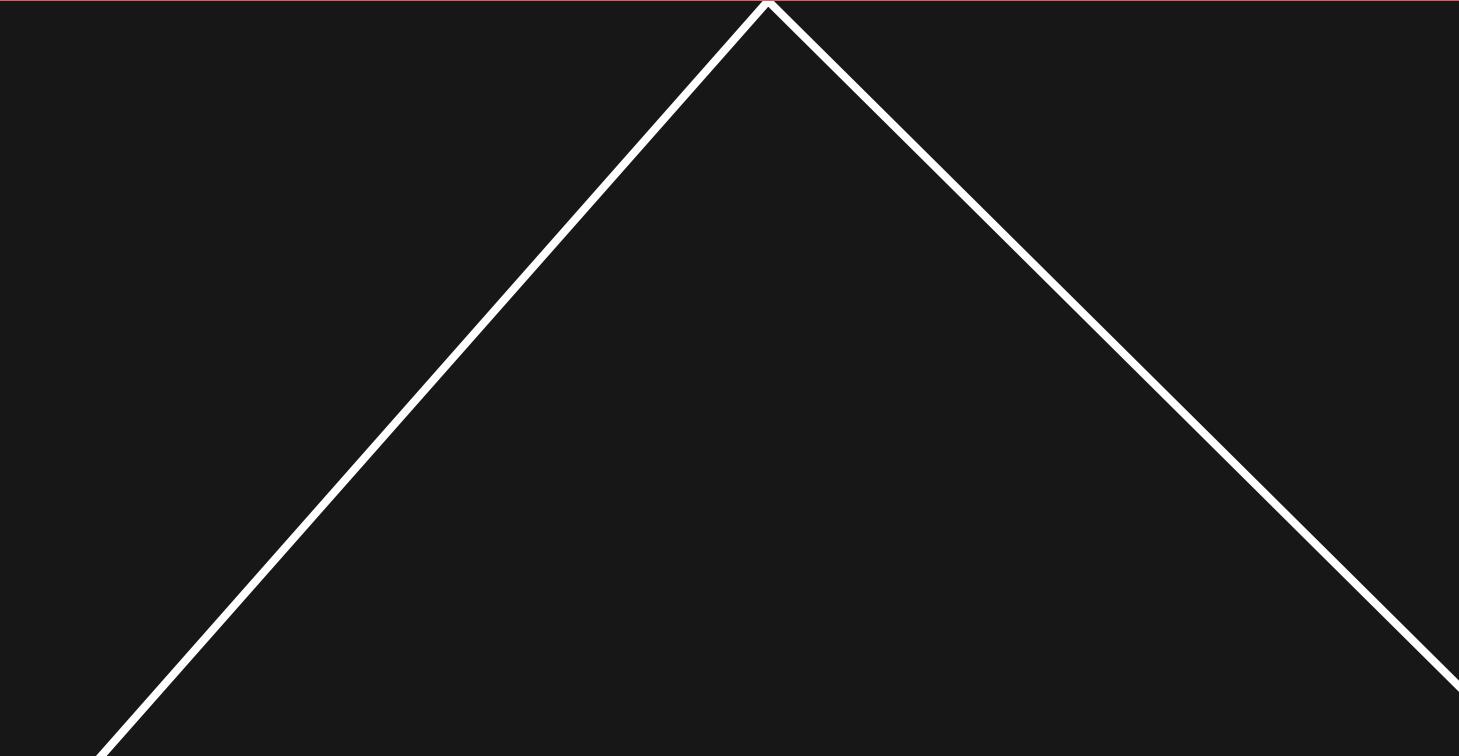
It is defined as variation in the DNA sequence leading to deletion, insertion, Duplication, substitution of nucleotides.



GENOMIC POLYMORPHISMS



GENOME WIDE
ASSOCIATION
STUDIES
(GWAS)

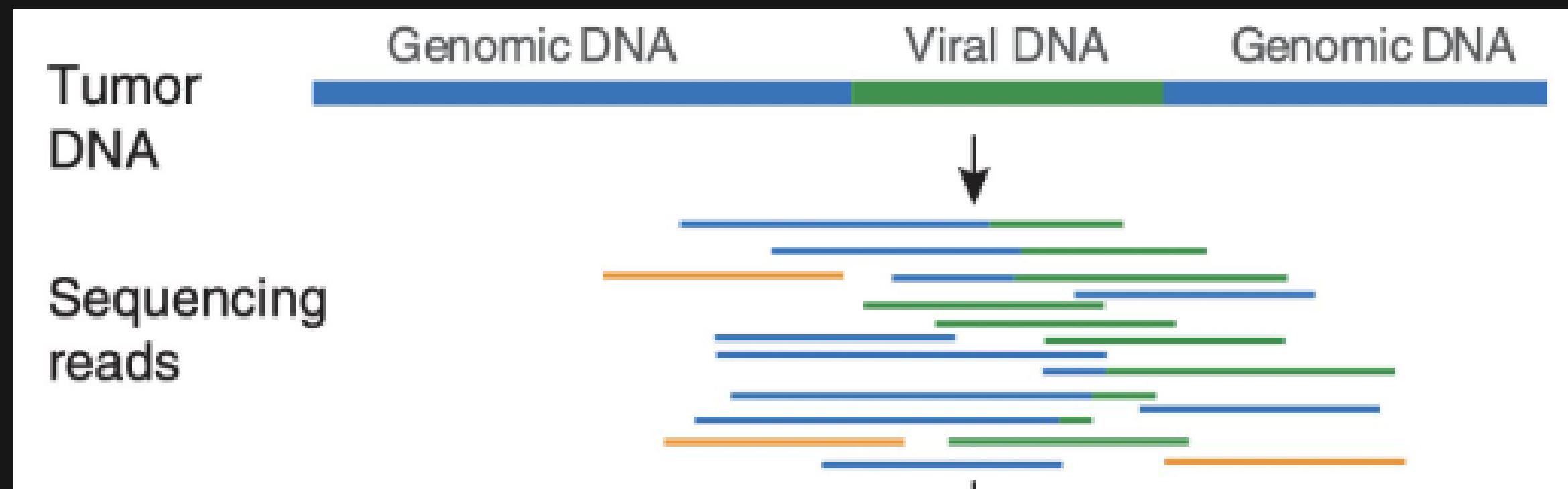


PHYLOGENETIC
TREE
RECONSTRUCTION

GWAS

- It is a research approach used to identify genomic variants that are statistically associated with a risk for a disease or a particular trait.
- More than 10% of reads were unmapped when mapping 12.2 million reads to the human genome using the popular Burrows-Wheeler Aligner.

- Some of the reads will be left unmapped due to errors generated during sequencing. Other reads are left unmapped for unknown reasons. It may be that some unmapped reads vary significantly
- Unmapped reads may contain more interesting and novel biological information than mapped reads because these reads diverge enough from the reference genome to remain unmapped



PHYLOGENETIC TREE RECONSTRUCTION

- Phylogenetic tree reconstruction is often completed through comparing homologous gene sequences in a group of species of interest.

- The method which we used to identify homologous genes is limited because it only allows for comparing gene sequences instead of comparing whole genomes.

METHODS



De Bruijn graph

**Relaxed
de Bruijn Graph**

DE BRUJIN GRAPH

- De Bruijn graphs are usually representative of a single species and are commonly used for genome assembly.
- In a de Bruijn graph, each node represents a unique kmer. Edges in the graph represent kmer overlaps. The graph is usually constructed from NGS reads where reads are broken into kmers and used to populate the graph .

Relaxed de Bruijn Graph

- The graph contains sequence information for multiple species.
- Kmers can occur multiple times in the graph.
- relaxing these constraints on the de Bruijn graph, we are able to identify interesting genomic variation in a tractable amount of time and space.

Algorithm:

De Bruijn graph:

Code:

```
def Construct(seq,k):
    for i in range(0,len(seq)-k+1):
        kmer = seq[i:i+k]
        kmer.append(kmer)
    return kmer

seq=input("Enter the dna sequencing")
k=int(input("Enter the k value"))
Construct(seq,k)

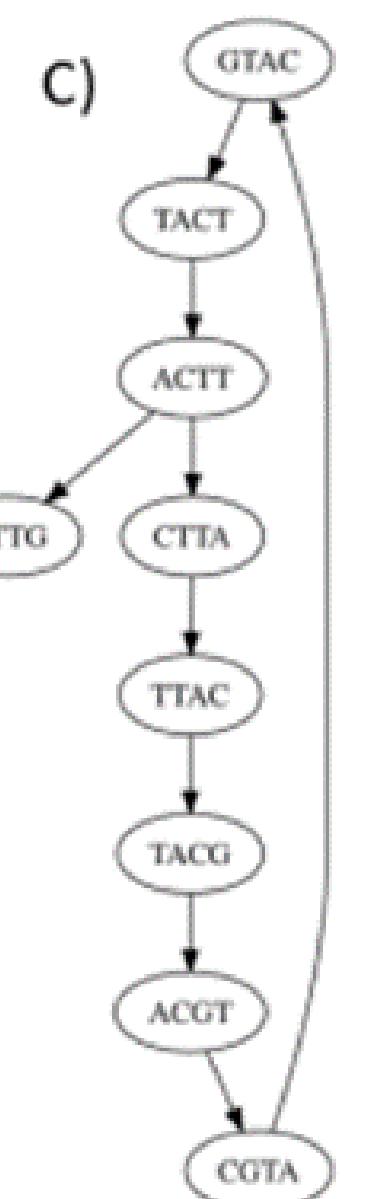
for i in kmer:
    su=i[-(k-1):]
    pre=i[:(k-1)]
    suffix.append(su)
    prefix.append(pre)
    graph.add_edge(su,pre)
```

Output:

A) **ACTTACGTACTTG**
ACTT
CTTA
TTAC
TACG
ACGT
CGTA
GTAC
TACT
ACTT
CTTG

B) -

B)



Relaxed de Bruijn Graph:

Algorithm 1:

Algorithm 1 Initial relaxed de Bruijn graph construction.

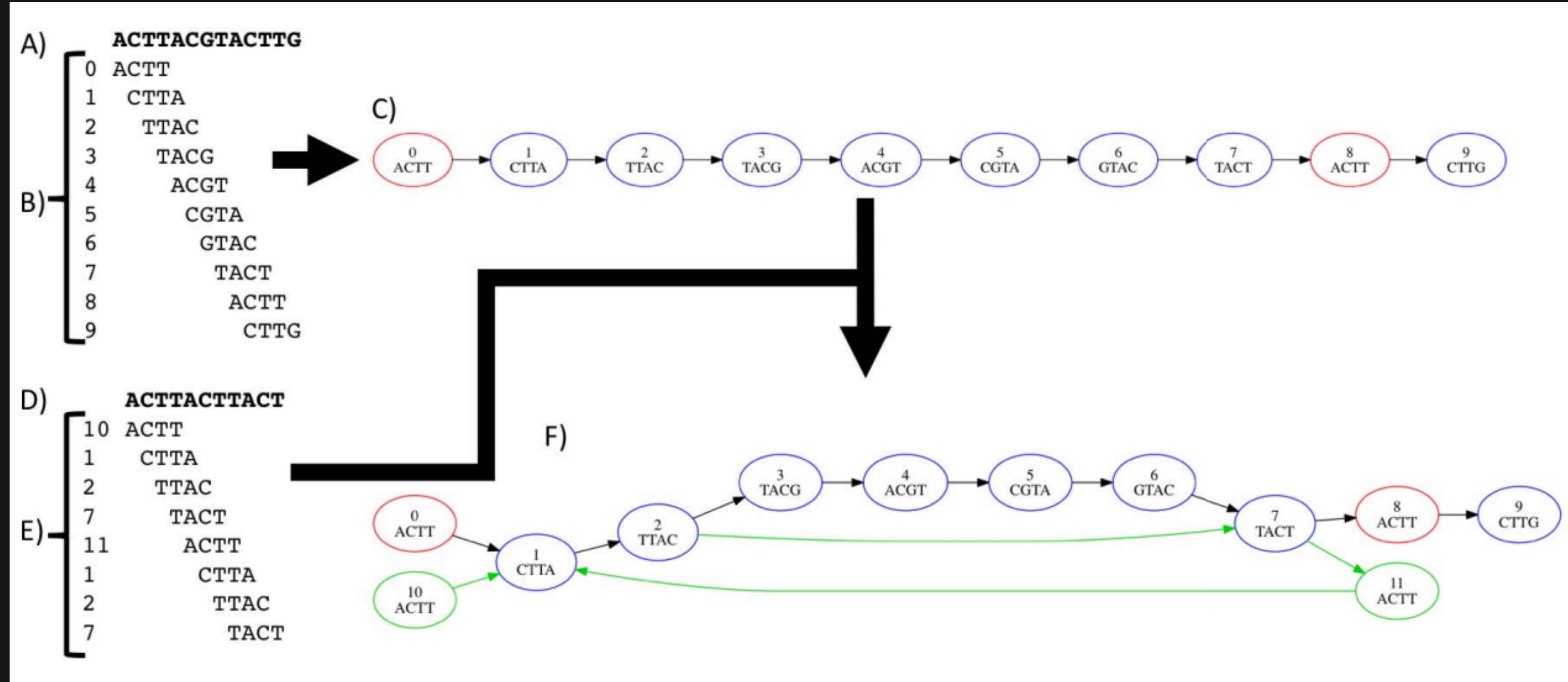
```
1: procedure CONSTRUCT(seq, k)
2:   Input: DNA sequence seq, kmer length k
3:   Output: kmer counts occs,
4:           index counter curidx,
5:           relaxed de Bruijn graph g
6:           kmer-index reverse lookup table rlookup
7:   occs  $\leftarrow$  occurrences of each kmer
8:   curidx  $\leftarrow$  0
9:   g  $\leftarrow$  an empty graph
10:  for each kmer kmer in seq do
11:    l  $\leftarrow$  prefix of kmer
12:    lidx  $\leftarrow$  curidx
13:    curidx  $\leftarrow$  curidx + 1
14:    occs[l]  $\leftarrow$  occs[l] + 1
15:    rlookup[l]  $\leftarrow$  lidx
16:    r  $\leftarrow$  suffix of kmer
17:    ridx  $\leftarrow$  curidx
18:    curidx  $\leftarrow$  curidx + 1
19:    occs[r]  $\leftarrow$  occs[r] + 1
20:    rlookup[r]  $\leftarrow$  ridx
21:    g.addedge(lidx, ridx)
22:  end for
23:  return occs, curidx, g
24: end procedure
```

Algorithm 2:

Algorithm 2 Appending new sequences after initial graph construction.

```
1: procedure APPEND(g, occs, curidx, seq, k, rlookup)
2:   Input: initialized relaxed de Bruijn graph g
3:           kmer occurrence counter occs
4:           index counter curidx
5:           DNA sequence seq
6:           kmer length k
7:           kmer-index reverse lookup table rlookup
8:   for each kmer kmer in seq do
9:     l  $\leftarrow$  prefix of kmer
10:    if occs[l] == 1 then
11:      lidx  $\leftarrow$  rlookup[l]
12:    else
13:      lidx  $\leftarrow$  curidx
14:      curidx  $\leftarrow$  curidx + 1
15:    end if
16:    r  $\leftarrow$  suffix of kmer
17:    if occs[r] == 1 then
18:      lidx  $\leftarrow$  rlookup[r]
19:    else
20:      ridx  $\leftarrow$  curidx
21:      curidx  $\leftarrow$  curidx + 1
22:    end if
23:    g.addedge(lidx, ridx)
24:  end for
25: end procedure
```

Observation:



RESULT:

- In the simplified graph, the mutation, insertion and inversion form simple bubble structures in the graph while the translocation forms a much more complex structure.
- In these very ideal conditions the generated graph shows structures that could be used to generate a phylogenetic signal or for phenotype association with additional generated graphs from other individuals.

CONCLUSION

- In this work, we have presented a method for constructing a unified, relaxed de Bruijn graph that contains more than one sequence source. The relaxed de Bruijn graph enables identification of graph structures that may be used as a signal for phylogenetic tree reconstruction or for use in association studies for phenotypes.

