



## Department of Electrical Engineering, Indian Institute of Technology Gandhinagar

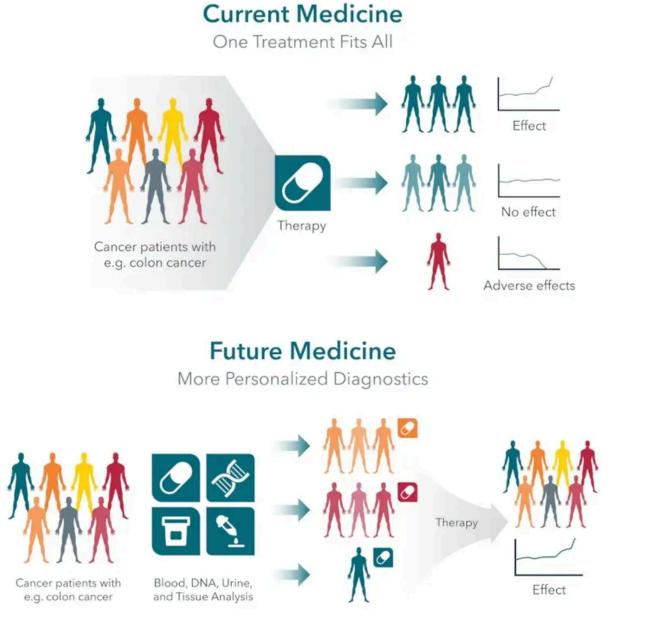
Anuja Chaudhari Sia Jariwala

**Ginisha Garg** Seemanshi Mall



### 1. Motivation

Genomics decodes the entire DNA sequence of an organism to uncover insights into genes, diseases, evolution, and personalized medicine.



#### **Challenges with Traditional Software-Based Genomic Analysis**

- Slow Processing: Software alone struggles with massive, evergrowing genomic datasets.
- High Computational Cost: Demands intensive CPU time and
- Data Privacy Concerns: Sensitive DNA data often sent to external cloud servers.
- Scalability Limits: Poor adaptability as sequencing throughput

### scales up.

#### **Practical Applications of Hardware-Accelerated Genomics**

Massively Parallel Processing: FPGAs handle

ideal for edge and embedded platforms.

of reads as they are sequenced.

reducing risk of external breaches.

multiple tasks simultaneously for ultra-fast analysis.

• Energy-Efficient Computing: Lowers power usage —

Real-Time Alignment: Enables immediate processing

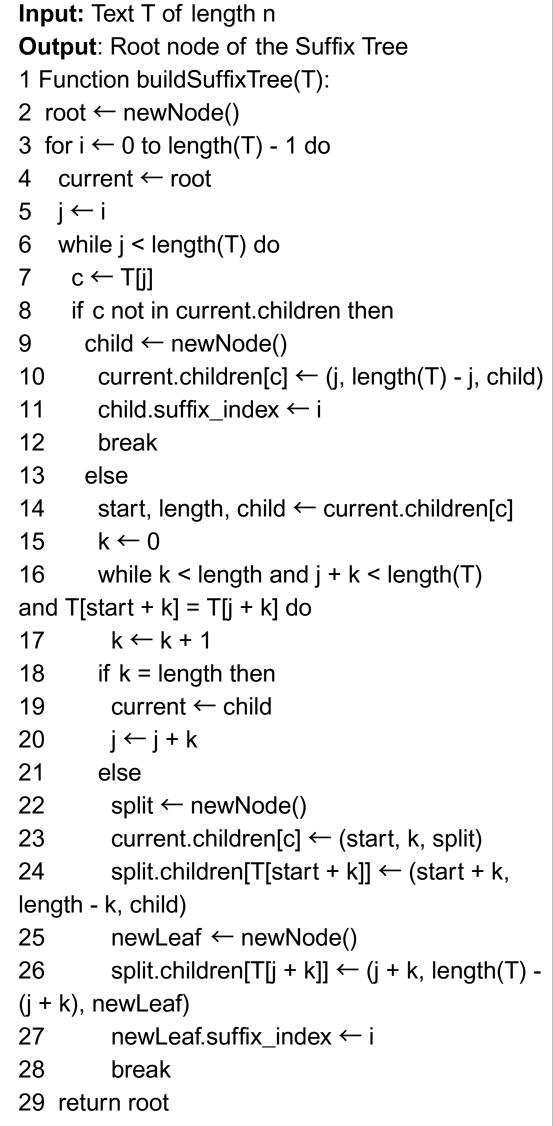
• Enhanced Privacy: Keeps genomic data on-device,

Why Hardware Acceleration?

**Build Suffix Tree from Text** 

- Clinical Diagnostics: Rapid identification of genetic disorders, cancer mutations, and pathogen genomes for real-time medical decisions. • Personalized Medicine: Tailoring treatment plans based on an individual's genomic profile — enabling more effective, targeted therapies.
- Pandemic Response: Fast sequencing and tracking of viral genomes (e.g., COVID-19 variants) to support vaccine and drug development.
- Ancestry & Evolutionary Studies: Large-scale comparison of genomes across populations/species to trace human ancestry and evolutionary patterns.
- Forensic Genomics: Rapid DNA matching for criminal investigations, missing persons identification, and disaster victim recovery.

## 2. Exact Matching



**Generate Suffix Array from Suffix Tree** 

if node.suffix index  $\neq$  -1 then

result.append(node.suffix\_index)

for key in sorted(node.children) do

\_, \_, child ← node.children[key]

Input: Root of Suffix Tree

1 Function getSuffixArray(root):

result ← emptyList()

Function DFS(node):

DFS(child)

Input: Text T, Suffix Array SA

BWT ← BWT + "\$"

 $BWT \leftarrow BWT + T[i - 1]$ 

**Compute BWT from Suffix Array** 

1 Function computeBWT(T, SA):

Result: Burrows-Wheeler Transform BWT

DFS(root)

10 return result

2 T ← T + "\$"

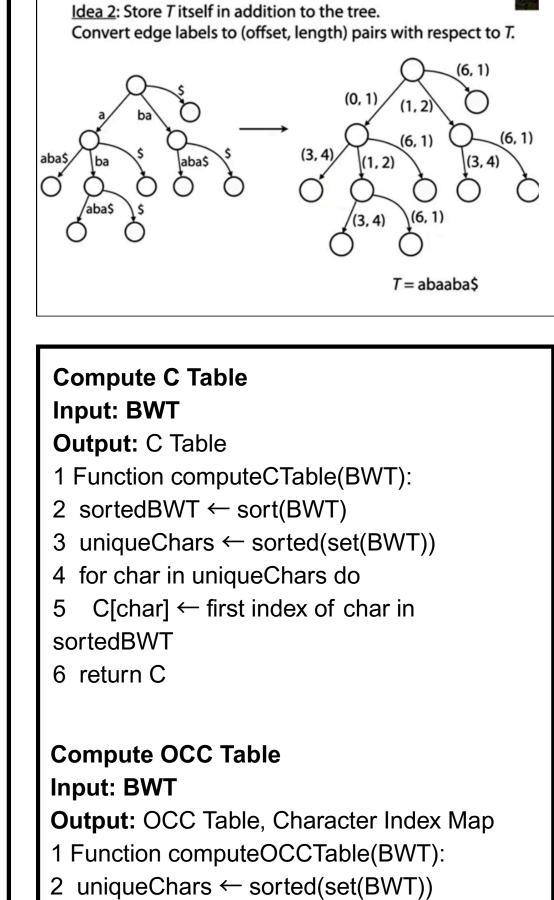
4 for i in SA do

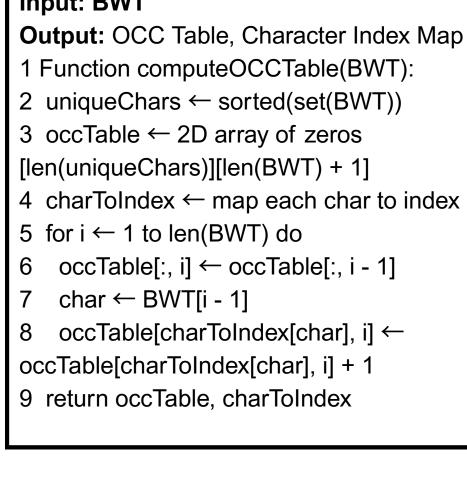
5 if i = 0 then

9 return BWT

3 BWT ← ""

**Output**: Suffix Array SA





0	banana\$	6	\$banan <mark>a</mark>
1	anana\$b	5	a \$ b a n a   n
2	nana\$ba	3	ana\$ba n
3	ana\$ban	1	a n a n a \$  <b>b</b>
4	na\$bana	0	b a n a n a  <b>\$</b>
5	a\$banan	4	na\$ban <mark>a</mark>
6	\$banana	2	nana\$b <mark>a</mark>



#### **Backward Search using FM-Index** Input: Pattern P, C Table, OCC Table, charToIndex, BWT Output: List of match indices in Suffix Array 1 Function backwardSearch (P, C, OCC, charToIndex, BWT): 2 s ← 1 $3 e \leftarrow len(BWT)$ 4 for char in reverse(P) do 5 if char not in charToIndex then return []

idx ← charToIndex[char] 8  $s \leftarrow C[char] + OCC[idx][s - 1] + 1$ 9  $e \leftarrow C[char] + OCC[idx][e]$ 10 if s > e then

return [] 12 return range(s - 1, e)

**Reconstruct Read from Seeds** 

#### Input: List of seeds, reference genome, FMindex data Output: Reconstructed sequence Function reconstructFromSeeds(seeds. refGenome, SA, C, OCC, charToIndex, BWT): reconstruction ← list of '-' of length(refGenome) 3 seedCount ← count each seed in seeds 4 for seed in seedCount: matches ← backwardSearch(seed, C OCC, charToIndex, BWT) positions ← [SA[i] for i in matches] used ← 0 for pos in positions do if used ≥ seedCount[seed] then

if pos + len(seed) ≤ length(refGenome) then reconstruction[pos:pos+len(seed)] ← list(seed) used ← used + 1 14 return join(reconstruction)

#### Applications of exact matching and fast backward search of query reads

**Mutation Detection:** Allows for quick identification of genetic mutations, enabling fast diagnostics and personalized treatments.

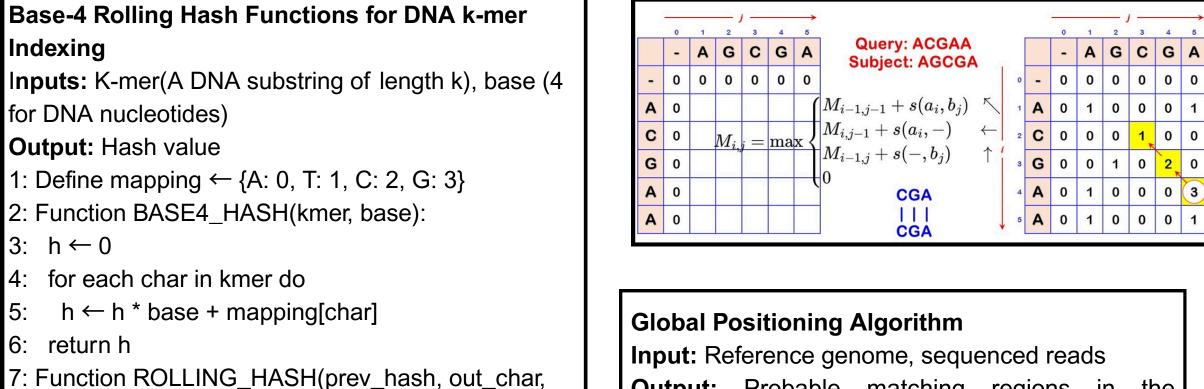
mutations .Quick detection of these mutations can help doctors choose targeted therapies and avoid ineffective treatments.

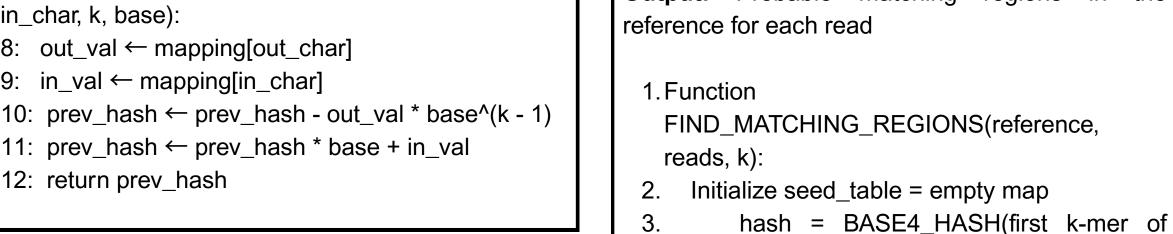
Some cancers or other diseases have specific

**Accelerating Gene Expression Analysis:** Quickly Mapping RNA Reads for Cancer Research and Identifying Potential Biomarkers

**CRISPR Off-Target Detection:** Ensures accurate genome editing by minimizing offtarget effects, enhancing the safety of gene therapies. It is widely used in crop improvement and in developing genetically modified animals for research.

# 3. Approximate Matching







Input: Reference Genome, reads and candidate position for each read Output: ref\_alignments: A map of reference positions

to aligned read bases. consensus\_sequence: A sequence of the most

frequent base at each position in ref\_alignments. Initialize ref\_alignments = empty map: reference position → list of aligned read bases

For each read and its candidate positions: Initialize best\_score = -∞

For each candidate position: Extract reference segment starting at that position

Align read with segment using SMITH\_WATERMAN Calculate alignment score 7. If score > best score:

8. Update best\_score 9. Store best\_alignment and best\_position

10. If best\_position is found: 11. Print aligned read and reference segment

12. For each aligned base pair in best\_alignment: 13. If reference base is not a gap:

14. Append read base to ref\_alignments at

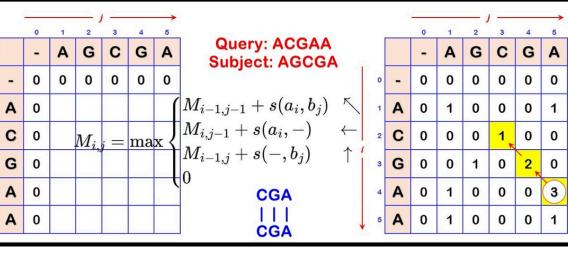
corresponding reference position # Generate consensus sequence

Initialize consensus\_sequence as empty string For each reference position in sorted(ref\_alignments):

If list of aligned bases is not empty: Find the most frequent base at that position Append it to consensus\_sequence

Else: Append 'N' to consensus\_sequence

reliability of genome analysis tools.



### Input: Reference genome, sequenced reads Output: Probable matching regions in the

reference) Store position 0 under hash

For i from 1 to len(reference) - k: hash = ROLLING HASH(hash, reference[i-1], reference[i+k-1], k)

Store position i under hash

For each read:

hash = BASE4\_HASH(first k-mer) Get all positions from seed\_table

matching the hash Return map of read → list of candidate positions

#### **Final Reconstruction**

**Input:** reference\_genome , ref\_alignments **Output:** Reconstructed Genome

1. For each position i from 0 to length of

Initialize reconstructed\_reference as empty string

reference\_genome - 1:

2. If ref\_alignments[i] is not empty: Count occurrences of each base in

ref\_alignments[i] Select the base with the highest frequency (most common base)

Append that base to reconstructed\_reference Else:

Append '-' to reconstructed\_reference to indicate missing coveragev

**Applications of Approximate Matching:** 

Error-Tolerant Read Mapping: Enables accurate alignment of DNA reads that contain sequencing errors or biological variations, improving the sensitivity and

### Genetic mutations A T G A C C A T G C C A T C C C ATACC

#### **Smith Waterman Algorithm**

Input: Two sequences to be aligned locally. Output: The optimal local alignment

between the two sequences Define match reward =+2, mismatch

penalty = -1, gap penalty = -1

• Initialize scoring matrix (4x4) for base alignment

Function GET\_SCORE(base1, base2): Return value from scoring\_matrix[base1][base2]

SMITH\_WATERMAN(seq1 seq2):

1. Initialize rows: previous\_row, current\_row (size = len(seq2) + 1)

2. Initialize max score=0 (0,0)max\_position= and traceback matrix  $(len(seq1)+1 \times len(seq2)+1)$ 

cell (i, j): Set H[i][j] = max(0, match/mismatch= H[i-1][j-1] score(A[i-1], B[j-1])

3. Fill Traceback matrix H: For each

,deletion= H[i-1][j] + gap\_penalty insertion= H[i][j-1] + gap\_penalty)

4. Start traceback from the cell with the maximum score, move:

Diagonally if the current cell was derived from a match/mismatch, Up if from a deletion, Left if from an insertion 6. Stop when a cell with value 0 is

reached and construct the aligned sequences during the traceback

# **Detection of Genetic Variants:** Facilitates identification of SNPs, insertions, deletions, and other structural variants, even when reads deviate slightly from the

reference genome — critical for population genetics and disease research. Robust Genome Assembly: Helps reconstruct genomes from fragmented or noisy reads by allowing flexible overlaps, ensuring more complete and accurate assemblies in de novo sequencing.

# 4. Flowchart

