

Abstract

Motivation

- Novel high-throughput sequencing methods
 - Huge sets of genomes with annotations
 - Opportunities for personalized genome-based medicine
- 1 genome \leq 1GB (1 base pair \approx 2 bits)
 - Pangenome of Dortmund (\approx 500.000 inhabitants) \approx \approx 500TB
 - Datasets too huge to maintain in commodity computers efficiently
- Genomes of two individuals are up to 99% identical [FCS06]
 - Opportunities for compression

Goal

Build an application, that

- maintains a pangenome
- answers queries like (approximate) pattern matching efficiently on the maintained pangenome

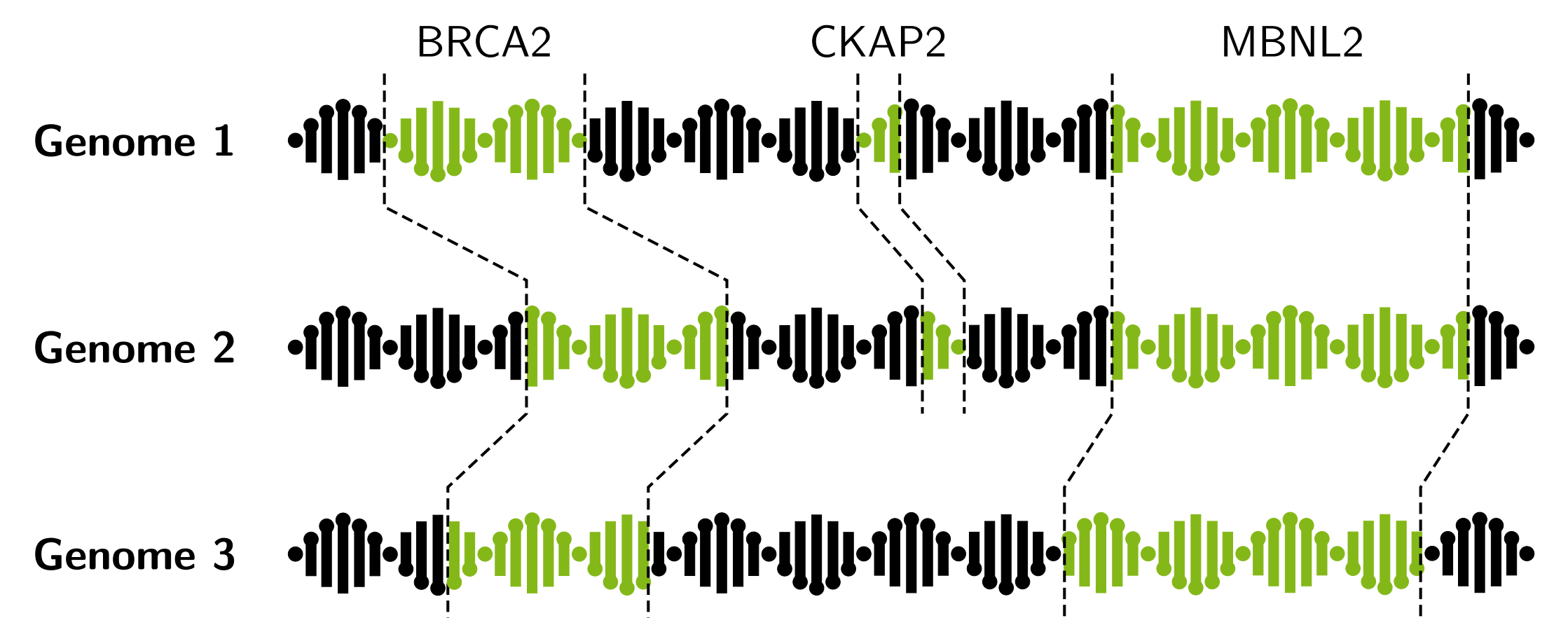
Implemented Data Structures

- Colored De Bruijn Graph [BBB+17] as a prefilter
- Journaled String Tree [RWR14]
- CHICO [Val16]
- Relative Lempel-Ziv [KPZ11]

Future Work

- Implement other approaches like grammar-based self-indices [CN12]
- Use annotations for annotation based queries or more efficient query answering

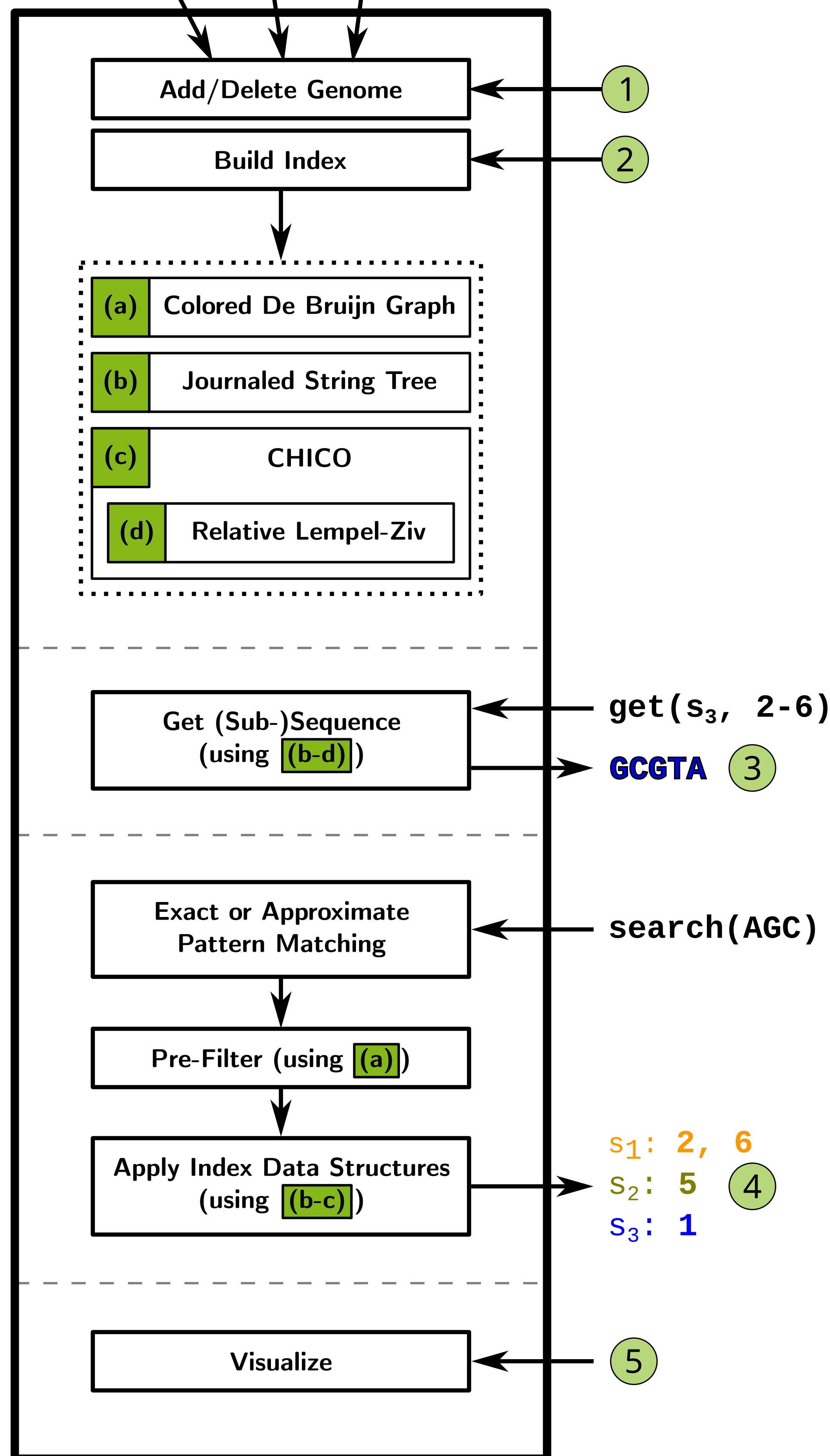
Annotations



Annotated area (for specific gene) is very similar among all genomes

- Highly compressible areas can be found with help of annotations
- Useful to fix positions for sequence alignment

Position: 1 2 3 4 5 6 7 8 9 10
 Genome s_1 : T A G C G T A G C A
 Genome s_2 : T A G G T A G C A
 Genome s_3 : T A G C G T A G T A



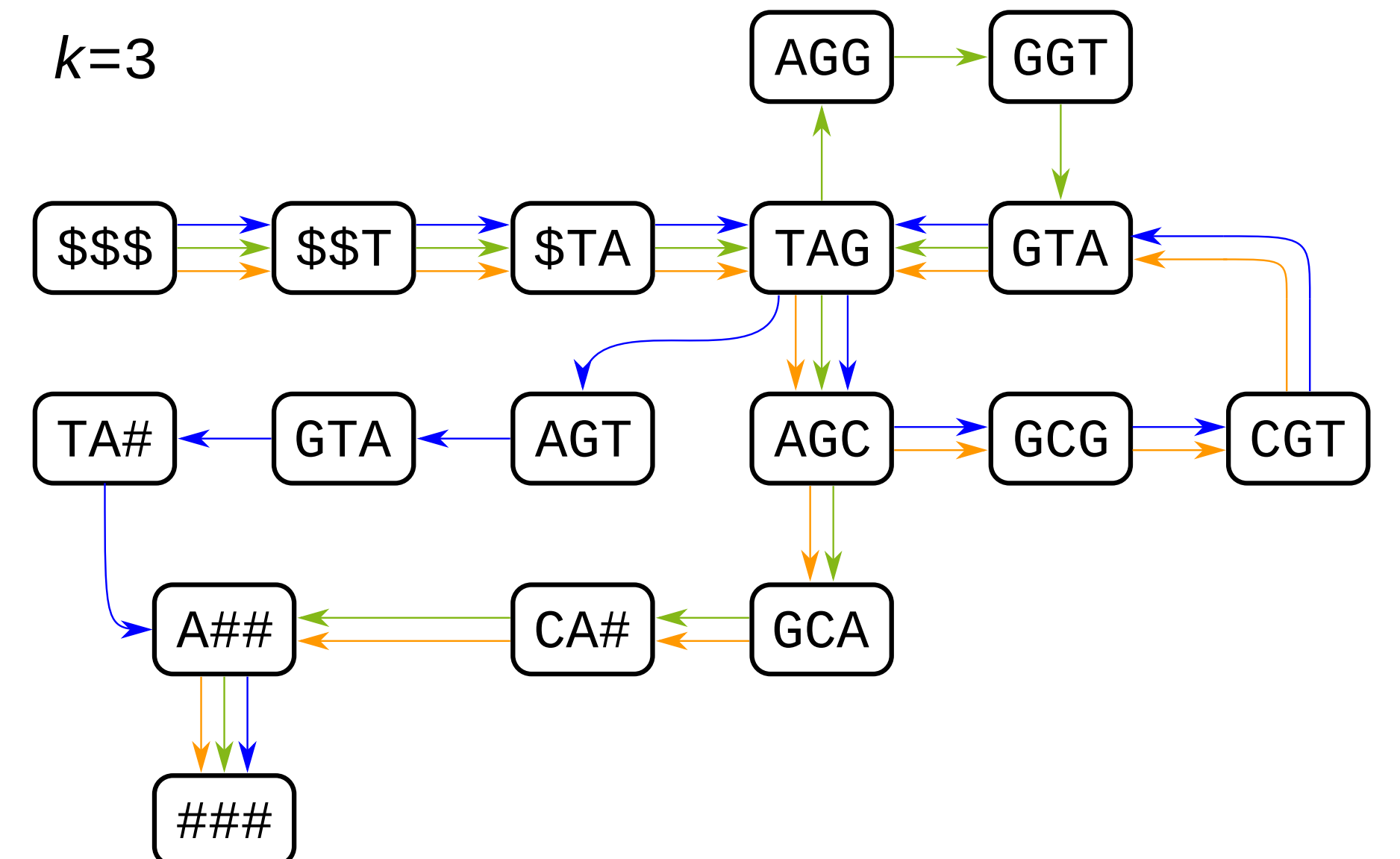
PanGeA Index

Index Features Usage

- 1 Add/Delete genomes
- 2 Build implemented data structures (a) - (d)
Existing structures are extended dynamically
- 3 Retrieve genomes or subsequences of genomes by ID and position or name of annotated area
- 4 Pattern Matching
Return a list of all occurrences of the given pattern in all genomes
- 5 Visualize the genome collection

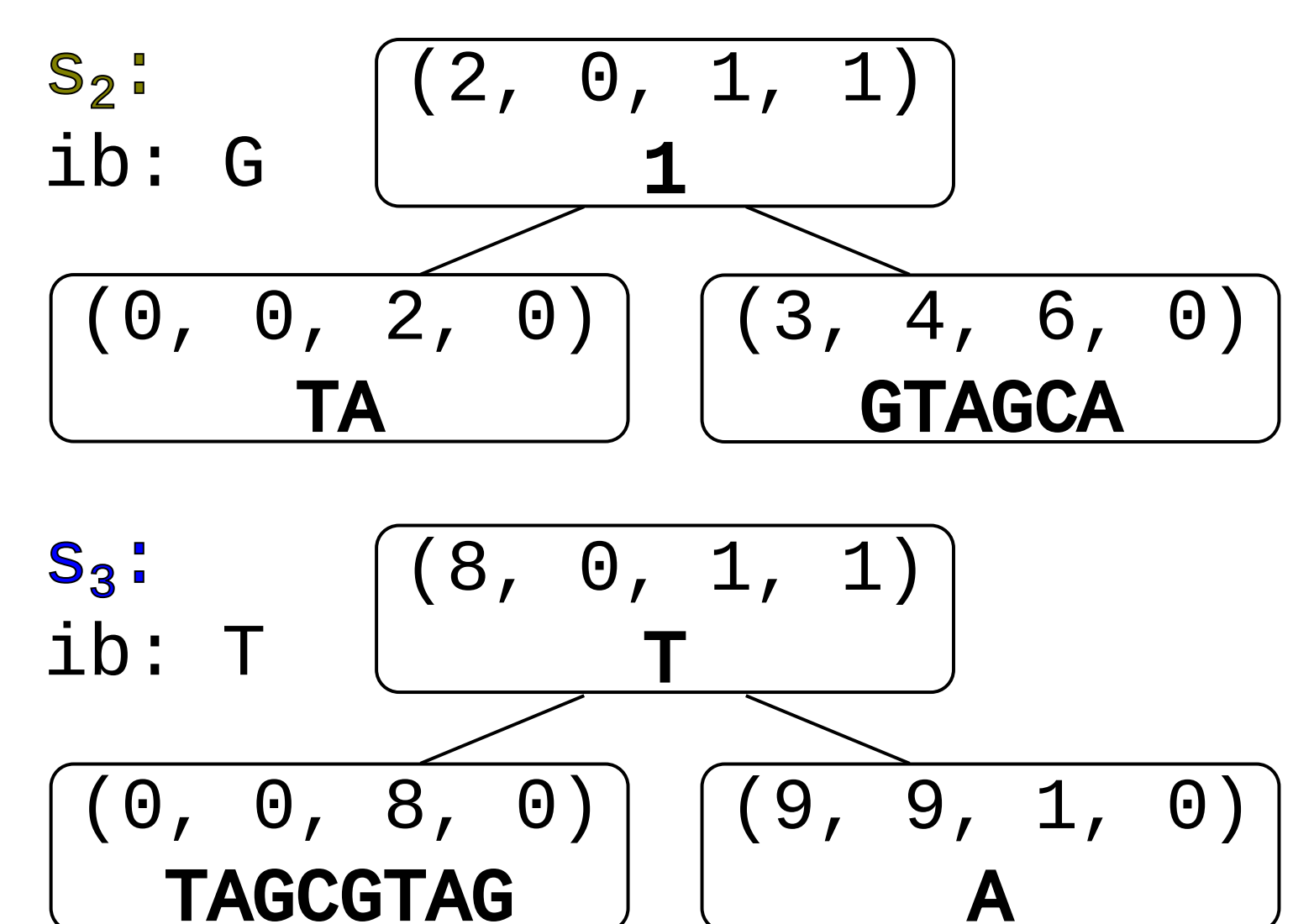
(a) Colored De Bruijn Graph:

- Edges represent $k+1$ -mers: $AGC \rightarrow GCG$
→ AGCG occurs in s_1 and s_3
- May return false positives due to circles
→ Use DBG as a pre-filter
- Supports approximate pattern matching
- Succinct representation allows memory efficient implementation [BOSS12]



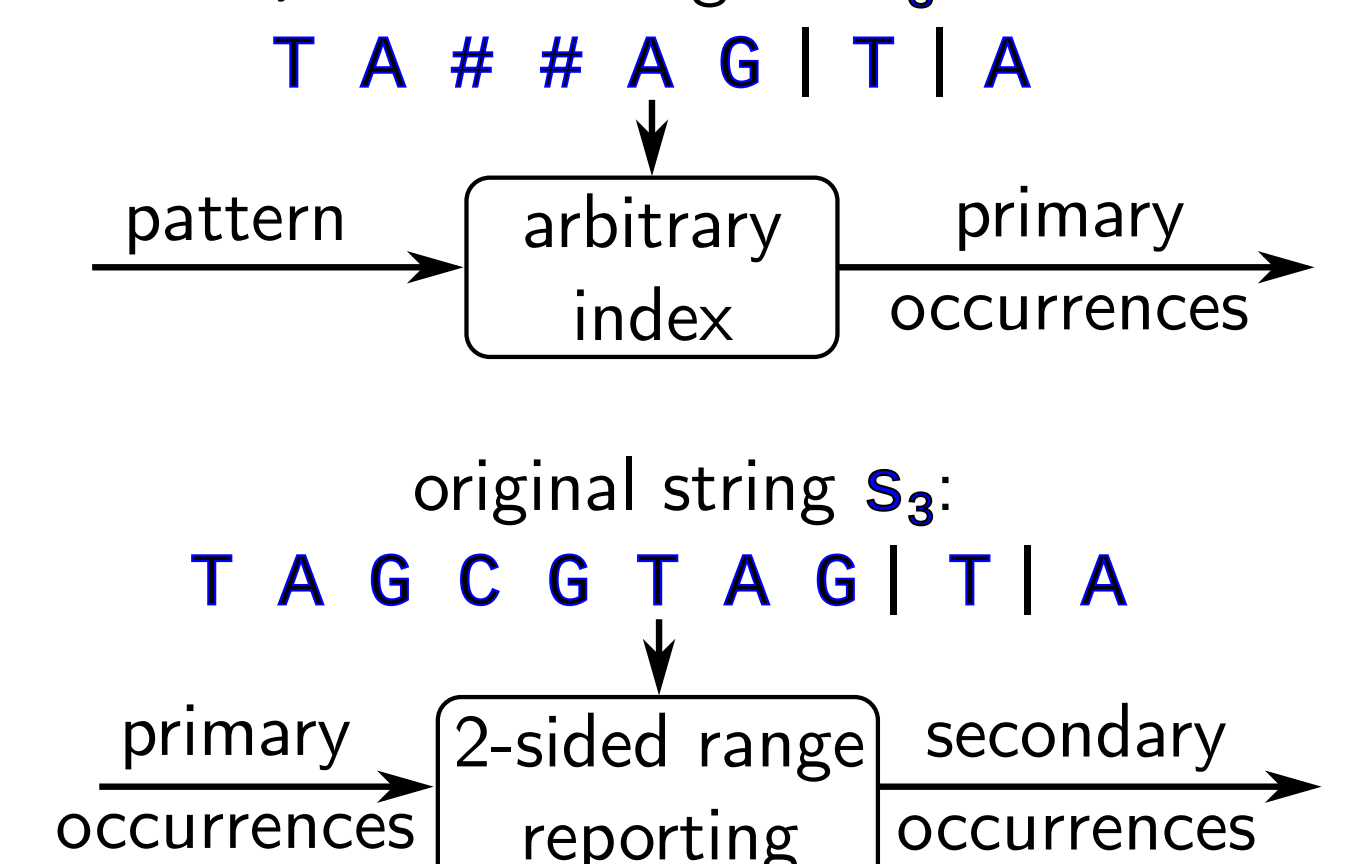
(b) Journaled String Trees

- Represent each genome as an order of copy operations that transform the reference genome into the other genome
- Save an insertion buffer for each genome and a tuple (vp, pp, l, t) for each copy operation
- Save tuples in balanced binary search tree



(c) CHICO (Compressed Hybrid Index for Repetitive Collections) kernel string for s_3 :

- Exact and approximate pattern matching
- Predefined maximum pattern length and maximum edit distance
- Utilizes RLZ-parsing to distinguish between two different possible match types
 - Primary occurrence: The pattern spans multiple RLZ-phrases
 - Secondary occurrence: The pattern is contained in a single RLZ-phrase



(d) Relative Lempel-Ziv

- User chooses reference genome r (here $r = s_1$)
- Each genome is interpreted as a composition of reference genome substrings
- Each substring is stored as a (start position, length)-tuple
- Retrieval of each genome in linear time to the genome length

