

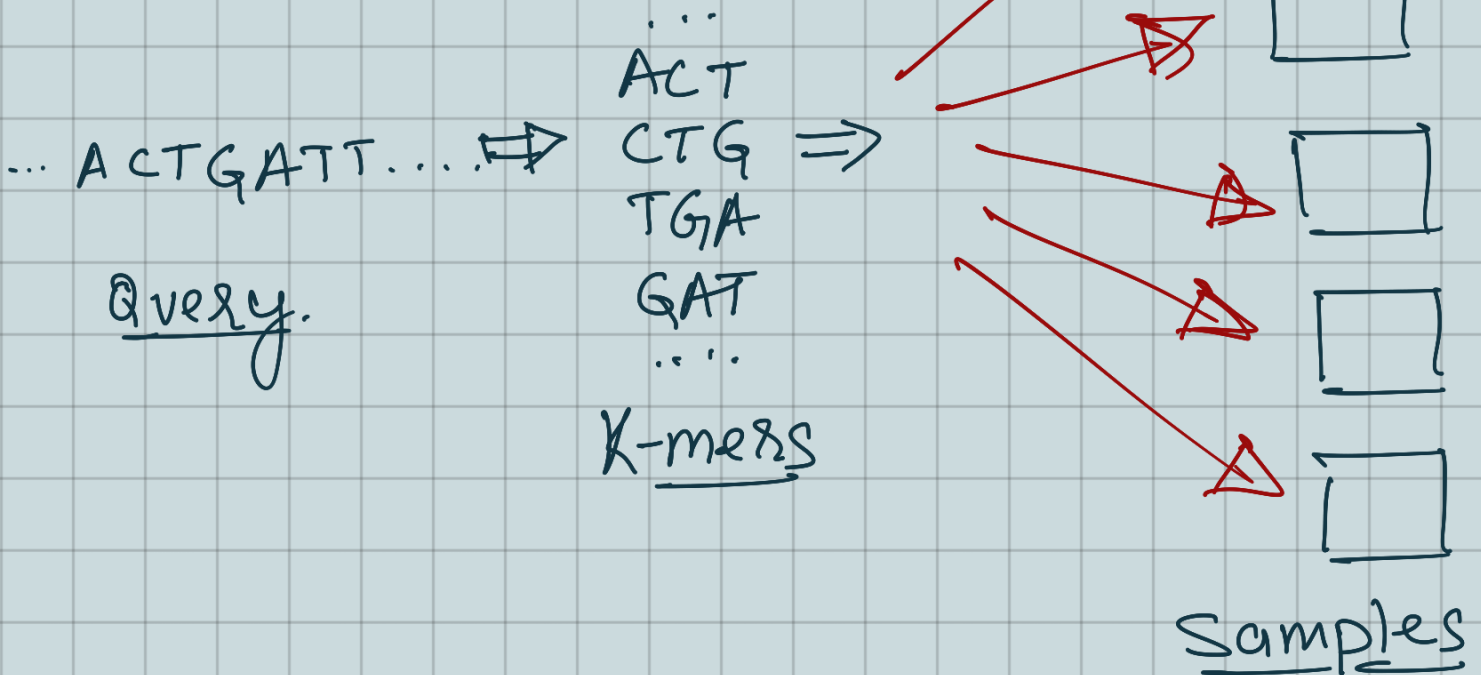
# Sample discovery problem:

SBT 2016

SSBT 2017

ALLSOME SBT 2017

MANTIS 2018



$\rightarrow$  Decompose each sample into k-mers

$\rightarrow$  If more than  $\theta$ -fraction k-mers from a query appear in a sample then there is a high chance that query appears in that sample.

$\rightarrow$  SBT, SSBT, ALLSOME SBT

$\rightarrow$  Based on Bloom filter.

$\rightarrow$  Bloom filter is used to represent k-mer contents.

# Mantis [Inverted index].

<u>S1</u>	<u>S2</u>	<u>S3</u>	<u>S4</u>
ACTT	AETG	AETG	
	TTTC	CTTG	CTTG
	GCGT	GCGT	GCGT
	AGCC	AGCC	

Map k-mers to samples.

→ Another layer of indirection



<u>k-mer.</u>	<u>Samples</u>
ACTG	S <sub>2</sub> S <sub>3</sub>
AETT	S <sub>1</sub>
CTTG	S <sub>2</sub> S <sub>4</sub>
TTTC	S <sub>2</sub> S <sub>3</sub>
GCGT	S <sub>2</sub> S <sub>3</sub> S <sub>4</sub>
AGCC	S <sub>2</sub> S <sub>3</sub> .

<u>k-mer.</u>
ACTG
AETT
CTTG
TTTC
GCGT
AGCC

<u>ID</u>
0
10
1
0
11
0



<u>IDs to Samples.</u>	
0	→ S <sub>2</sub> S <sub>3</sub>
1	→ S <sub>3</sub> S <sub>4</sub>
10	→ S <sub>1</sub>
11	→ S <sub>2</sub> S <sub>3</sub> S <sub>4</sub> .



	S <sub>1</sub>	S <sub>2</sub>	S <sub>3</sub>	S <sub>4</sub>	
0	→	0	1	1	0
1	→	0	0	1	1
10	→	1	0	0	0
11	→	0	1	1	1

## Scalability Challenge:

→ Mantis index provides fast query, and scales well up to thousands of input experiments/samples.

→ But we really want to index on the order of  $10^5 - 10^6$  samples.

### Key observation:

→ K-mers grow at worst linearly.

→ Color classes grow super-linearly.

★ Need a fundamentally better color class encoding.

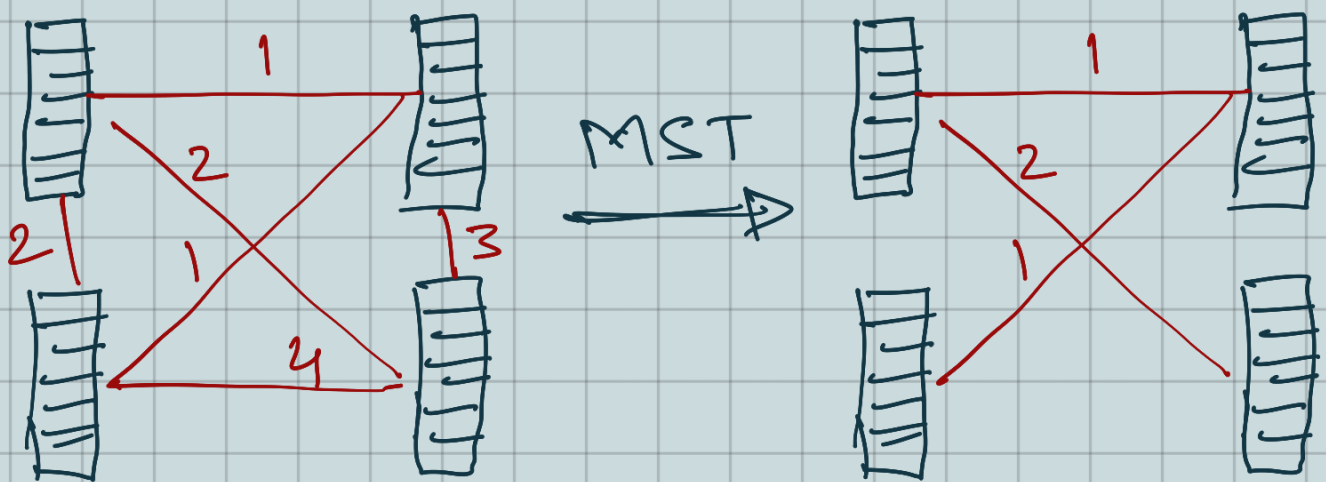
### Key idea:-

→ Exploit coherence between rows of the color class matrix

→ If we know a k-mer's color vector, we probably know a lot about its neighbor's color vectors.

## Build a color class Graph:

- Each color class is a vertex.
- Every pair of color classes is connected by an edge whose weight is the hamming distance between the color class vectors.



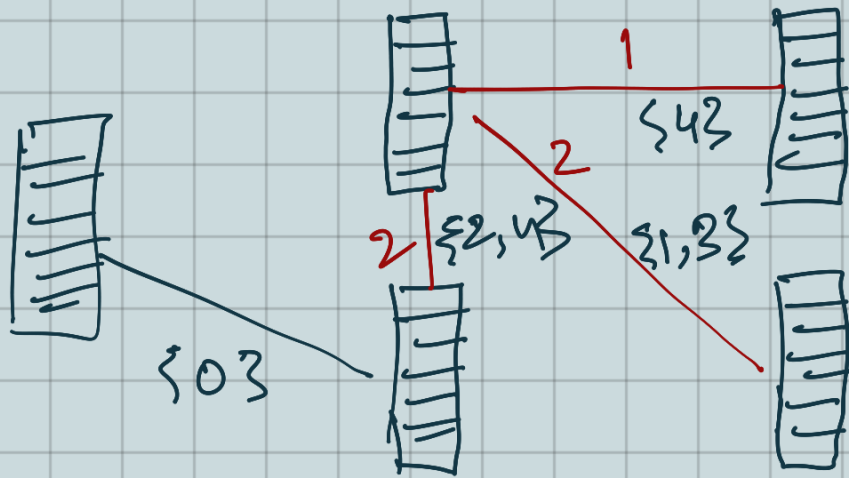
### Unfortunately:-

- There are many color classes (full graph too big)
- They are high dimensional (millions of samples)
- Neighbor search is very hard (LSH scheme seem to work poorly).

## Implicitly Represent the colored dBG.

- k-mers form the dBG.
- Can explicitly query neighbors.
- Use deBruijn Graph as an efficient guide for near-neighbor search in the space of color classes
- dBG helps to impose biological relationships over color-class vectors.
- MST derived from dBG-based color class graph is very close to the MST from complete graph.

MST efficiently encodes related color classes



→ Along the edge from each  $c$  node to its parent  $p$ , store  $S(p,c)$

→ the positions of the bits whose parity we would have to flip to obtain the child from the parent.

→ To reconstruct a node's bit vector, walk from the node to the root, flipping the parity of the positions you encounter on each edge