

Ex 1

 $S = \{\text{bio, bionic, bit, litly, buzz, car, caso, cast, zoo}\}$

The main idea of algorithm: divide alphabetically ordered strings in groups of size which can be stored inside the memory, then take the first strings of each group-block and create a ~~tree~~ prefix-tree. So when we will need to find our query we would be able to find the necessary block using the tree.

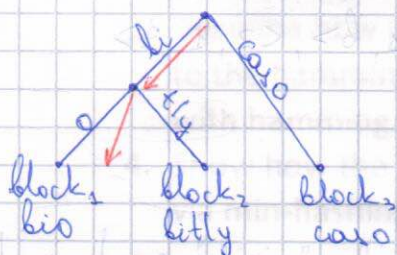
1) We divide our set in 3 blocks. We can compress the data inside the blocks so they will take less space on the disk:

 $\text{block}_1 = \{\text{bio, bionic, bit}\} = \{\text{bio, 3nic, 2t}\}$ (so 3nic is equal to

take 3 first letters of last word and add nic; this compression will give good results because words are alphabetically ordered)

 $\text{block}_2 = \{\text{litly, buzz, car}\} = \{\text{litly, 1uzz, 0car}\}$
 $\text{block}_3 = \{\text{caso, cast, zoo}\} = \{\text{caso, 3t, 0zoo}\}$

2) Now we build a tree. Let's assume that the number of blocks isn't big so we can store in memory the tree based on first element of each block



Execution of search:

With our assumption we can expect that we will get result with going to disk only once (to get the block where we might find our string).

Search of string $P = \text{bis}$:

1) The track of searching in the tree is shown with the red pen. So we will first go to branch bi and then our string must be somewhere between two branches because $o < s < t$, so we must search the P inside the block_1

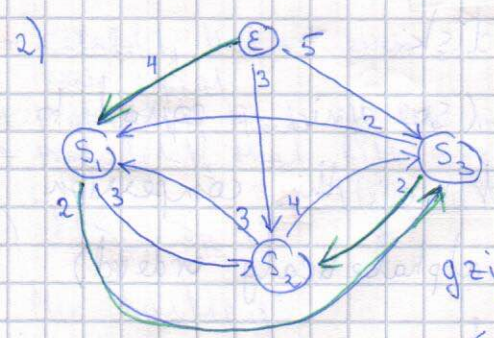
2) From disk we get $\text{block}_1 = \{\text{bio, 3nic, 2t}\}$. Scanning through the block we decompress it and check if any string is equal to our searched P. In this case we won't find one. This means there is no string P in set S.

Ex 2.

$s_1 = \text{"abaco"}$ $s_2 = \text{"dada"}$ $s_3 = \text{"coaba"}$

1) $gzip(s_1 s_2 s_3) = gzip(\text{abacodadacoaba}) = \langle 0, 0, a \rangle;$
 $w=3$ $\langle 0, 0, b \rangle;$
 $\langle 2, 1, c \rangle;$
 $\langle 0, 0, 0 \rangle;$
 $\langle 0, 0, d \rangle;$
 $\langle 0, 0, a \rangle;$
 $\langle 2, 2, c \rangle;$
 $\langle 0, 0, 0 \rangle;$
 $\langle 3, 1, b \rangle;$
 $\langle 2, 1, \$ \rangle;$

↑
 symbol of end of string, because we don't need to add anything else.



$gzip(s_1) = gzip(\text{abaco}) = \langle 0, 0, a \rangle; \langle 0, 0, b \rangle; \langle 2, 1, c \rangle;$
 $\langle 0, 0, 0 \rangle$

$gzip(s_2) = gzip(\text{dada}) = \langle 0, 0, d \rangle; \langle 0, 0, a \rangle; \langle 2, 2, \$ \rangle$

$gzip(s_3) = gzip(\text{coaba}) = \langle 0, 0, c \rangle; \langle 0, 0, 0 \rangle; \langle 0, 0, a \rangle;$
 $\langle 0, 0, b \rangle; \langle 2, 1, \$ \rangle$

$gzip(s_1/s_2) = gzip(\text{abaco/dada}) = \langle 0, 0, d \rangle; \langle 4, 1, d \rangle; \langle 2, 1, \$ \rangle$

$gzip(s_1/s_3) = gzip(\text{abaco/coaba}) = \langle 2, 2, a \rangle; \langle 7, 2, \$ \rangle$

$gzip(s_2/s_1) = gzip(\text{dada/abaco}) = \langle 1, 1, b \rangle; \langle 2, 1, c \rangle; \langle 0, 0, 0 \rangle$

$gzip(s_2/s_3) = gzip(\text{dada/coaba}) = \langle 0, 0, c \rangle; \langle 0, 0, 0 \rangle; \langle 3, 1, b \rangle; \langle 2, 1, \$ \rangle$

$gzip(s_3/s_1) = gzip(\text{coaba/abaco}) = \langle 3, 3, c \rangle; \langle 8, 1, \$ \rangle$

$gzip(s_3/s_2) = gzip(\text{coaba/dada}) = \langle 0, 0, d \rangle; \langle 2, 3, \$ \rangle$

So now we should find the "trace" in graph which will "visit" each node and will weight the less. In this case we have two routes with equal minimum weight 8, I will take randomly one: s_1, s_3, s_2 .

So Δ compression will be like this: $gzip(s_1), gzip(s_1/s_3), gzip(s_3/s_2)$