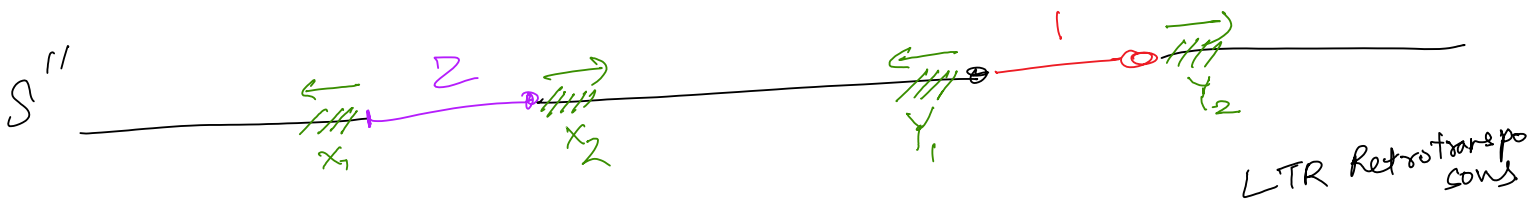
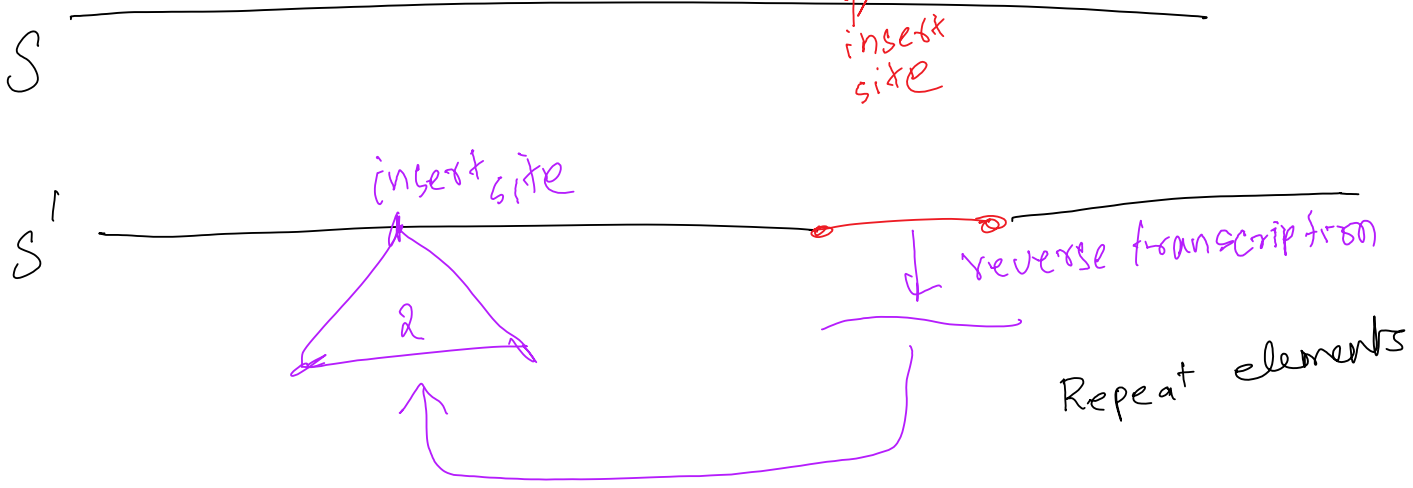
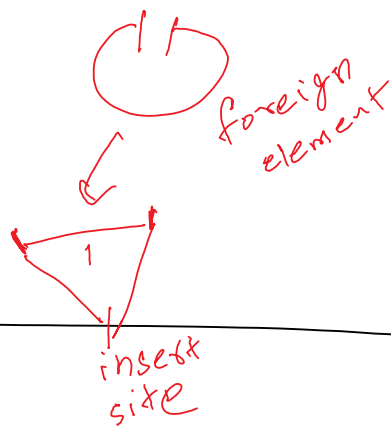
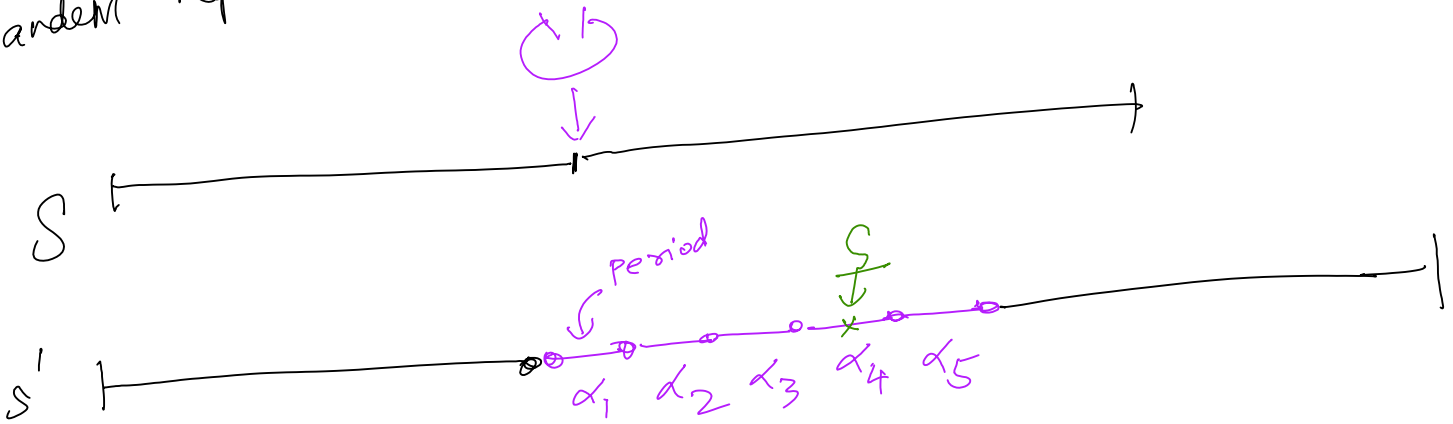


# Exact Matching

Monday, February 22, 2021 10:22 AM



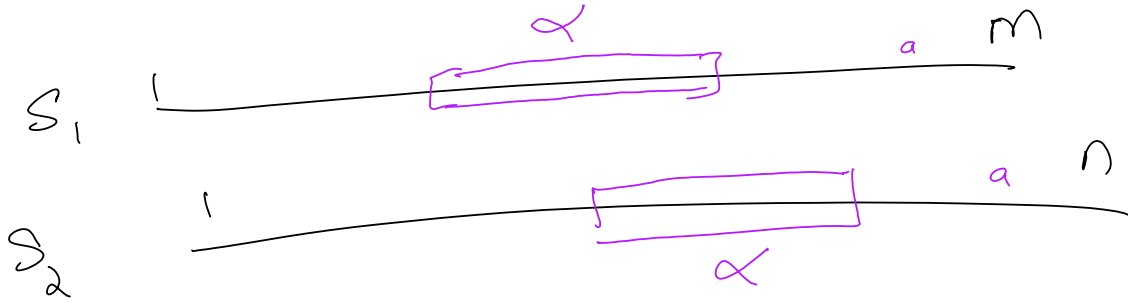
Tandem Repeats:



# Exact matches between $\geq 2$ sequences

Monday, February 22, 2021 10:23 AM

Use case #2: Multiple ( $\geq 2$ ) sequences (input)



Q) Find the longest common substring between  $S_1$  &  $S_2$ :

DP approach:

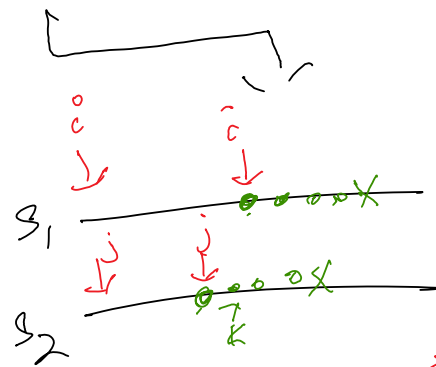
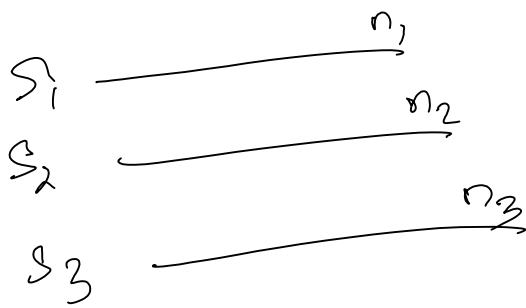
local  
 $\theta$ : parameters  
 $m_a = +1$   
 $m_i = -\infty$   
 $h = -\infty$   
 $g = -\infty$

LCS ( $S_1, S_2$ ) // exact matching reduced to

local Align ( $S_1, S_2$ )  
 $w/ \theta \{+1, -\infty, -\infty\}$

// inexact matching.  
 $O(mn)$  time  $n^2$

$O(n_1 n_2 n_3)$



for (i = ...  $n$ )  
 for (j = ...  $n$ )  
 for (k = ...  $n$ )  
 ++ matchy

use k-mers as building blocks to identify  $\alpha$

$O(mn \times n)$   
 $n^3$

Q) Can lcs be solved in  $O(m+n)$  time?

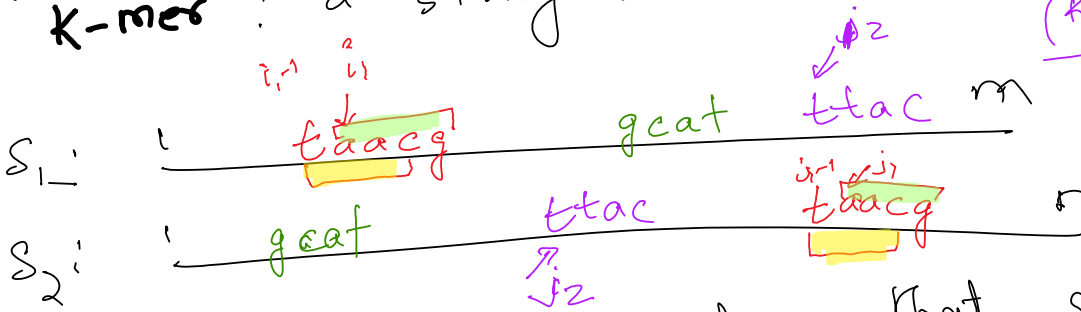


# Fixed length matches (k-mers)

Monday, February 22, 2021 10:23 AM

Finding exact matches of a Fixed length.  $k > \emptyset$

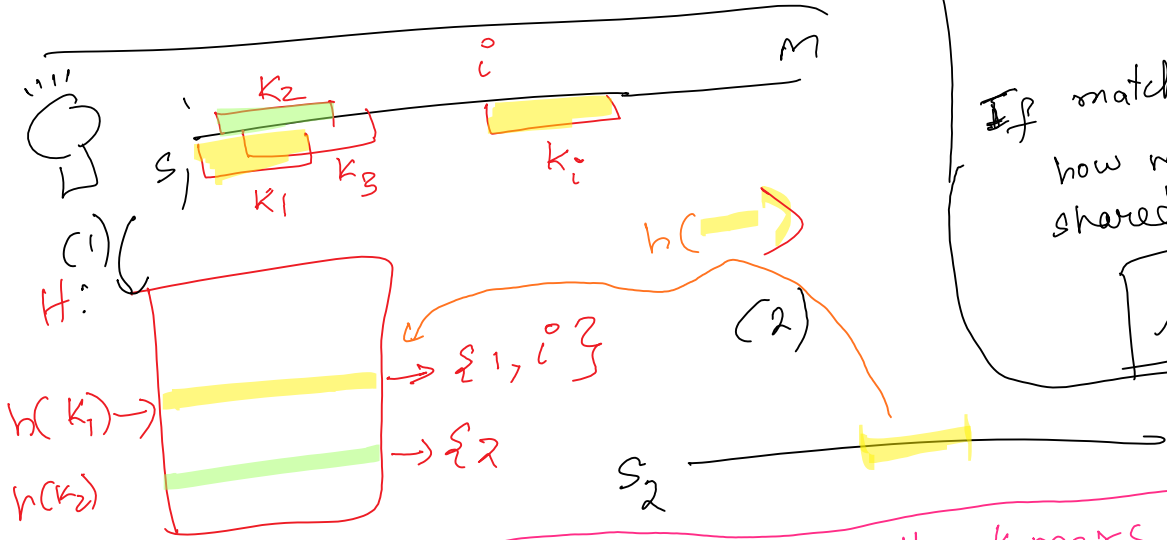
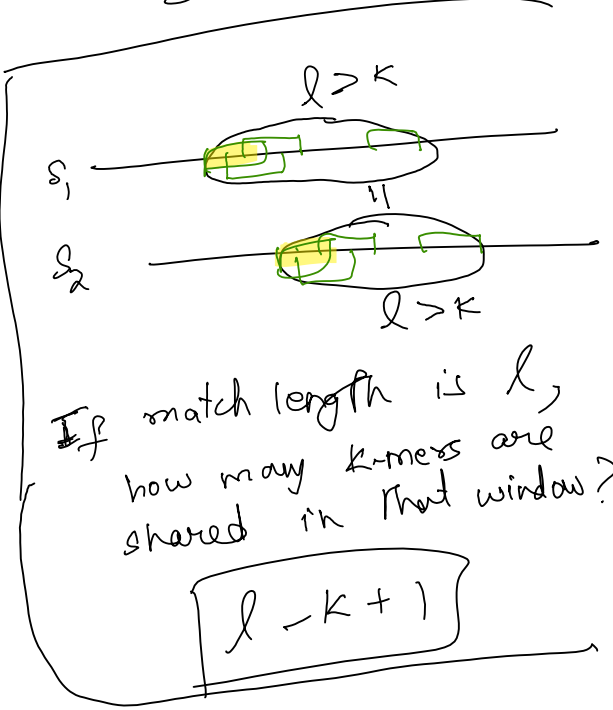
"k-mer": a string (or a substring) of length  $k > \emptyset$



n-grams  
q-grams

Goal: to identify all locations that start with the same k-mer

- $\langle S_1: i_1, S_2: j_1 \rangle$
- $\langle S_1: i_2, S_2: j_2 \rangle$
- $\vdots$
- $\langle S_1: i_{l-1}, S_2: j_{l-1} \rangle$



- Build a hash table/index for all k-mers in  $S_1$
- For each k-mer in  $S_2$ , search in  $H$  and retrieve the list of  $S_1$  indices.

# Look-up Tables

Wednesday, February 24, 2021

10:23 AM

## lookup tables

A collision-free hash table for all  $k$ -mers  
 $\Sigma = \{a, c, g, t\}$

$k=2$

$\Sigma = \{a, c, g, t\}$

$S: a c c g t a g$

$h("ac") = 1$        $h("gt") = 11$

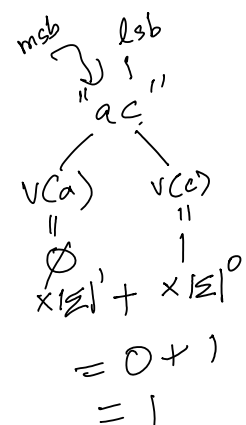
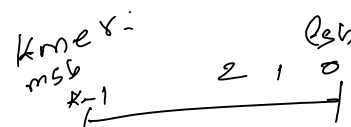
H:

0	aa
1	ac
2	ag
3	at
4	ca
5	cc
6	cg
7	ct
8	ga
9	gc
10	gg
11	gt
$ \Sigma ^k$	tt

$$2 \times |\Sigma|^1 + 3 \times |\Sigma|^0 = 11$$

lookup table size  $\neq$

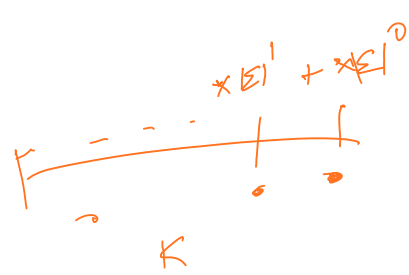
$$= |\Sigma|^k$$



$\Sigma = \{a, c, g, t\}$

00	= 0
01	= 1
10	= 2
11	= 3

symbol $c$	value $v(c)$
a	0
c	1
g	2
t	3



# Look-up table construction

Wednesday, February 24, 2021 11:44 AM

Q) How to build lookup table for the database string  $D$   
 $D = \{s_1, s_2, \dots, s_n\}$

Input:  $D[0 \dots N-1], K, \Sigma$

Output:  $LT(D)$

Algo:

1) Init.  $LT$  with  $|\Sigma|^K$  locations  $O(|\Sigma|^K)$

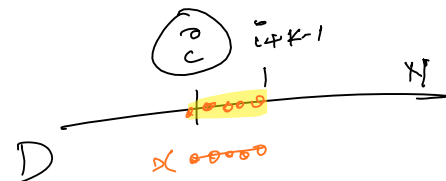
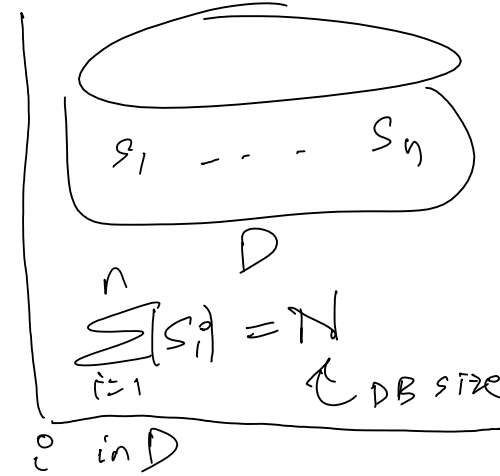
2) for ( $i = 0$  to  $N - K + 1$ )  $\rightarrow O(N)$

$x \leftarrow$  extract  $K$ -mer starting at  $i$

$\text{insert}(x, LT) \rightarrow O(K)$

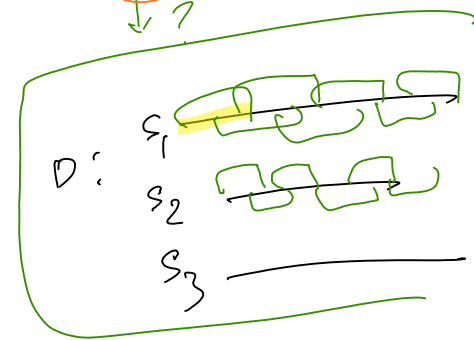
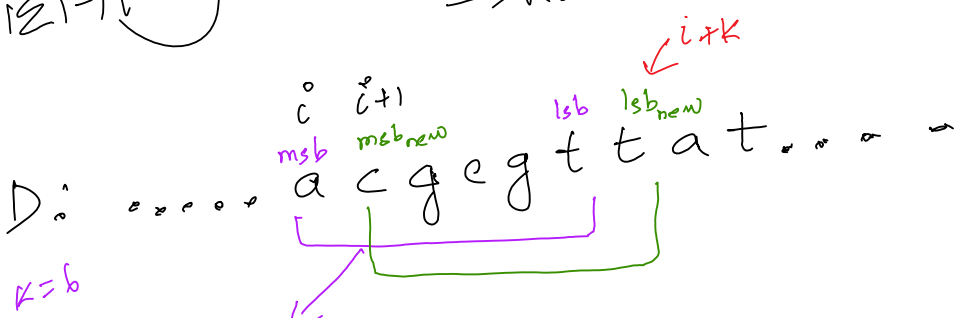
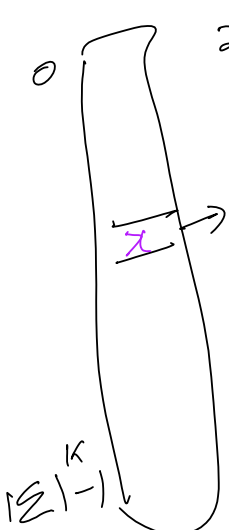
$h(x) \rightarrow$  index

$LT[\text{index}].\text{append}(i)$



$K = 11$   
 $= 12$

$\Rightarrow \text{Runtime} = O(|\Sigma|^K + N \cdot K)$



$x \leftarrow h("acgegt")$

$h_{\text{new}} \leftarrow \left( h_{\text{prev}} - \text{value}(D[i]) \times |\Sigma|^{K-1} \right) \times |\Sigma| + \text{value}(D[i+k])$

$\Rightarrow$  can calculate new hash value (at  $i$ ) in  $O(1)$  time (using the previous value at  $i-1$ )

$\Rightarrow$  construction time  $= O(|\Sigma|^K + N)$  time.

# Look-up table: Analysis

Friday, February 26, 2021 10:49 AM

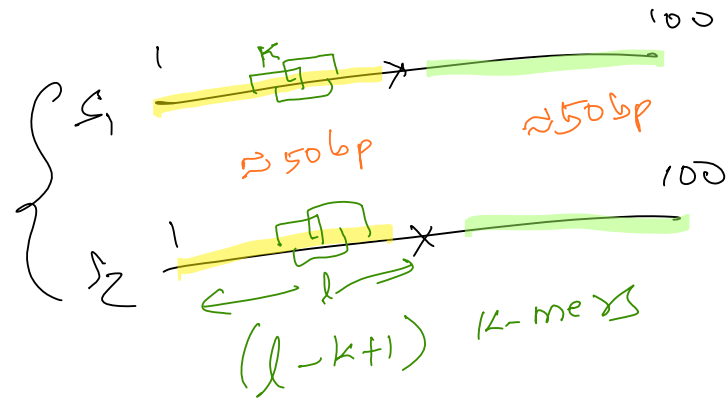
DNA:  $\Sigma = \{a, c, g, t\}$

k:	5	10	11	15	20
Size:	$4^5$	$4^{10}$	$4^{11}$	$4^{15}$	$4^{20}$
(#entries)	$\approx 2^{10}$	$\approx 2^{20}$	$\approx 2^{22}$	$\approx 2^{30}$	$\approx 1B$
	$\approx 1K$	$\approx 1M$	$\approx 4M$		

high expense  $|\Sigma|^k$

(k=11 used by BLAST)

reads



Q) variable length in a direct ways.