

In *genetics*, a sequence motif is a pattern that occurs in multiple sequences. Motifs are associated with important biological functions and finding “approximate” motifs is an important problem in Computer Science. In this assignment, you will implement one version of motif finding problem across various sequences.

Consider a set of n l -bit integers ($l \leq 64$): $\{S_1, S_2 \dots S_n\}$. Given a number d ($d \leq l$), the goal is to find all integers which are at a Hamming distance at most d from all of the given n integers. Hamming distance between two sequences (in our case, bit representation of integers) of equal lengths is the number of corresponding positions where the sequences differ.

A brute force approach for the problem would be to enumerate all possible 2^l integers and for each of these, check if it is within Hamming distance d from each of the n integers in the input set. For $l = 40$, this search space consists of more than 1 trillion entries, making it computationally infeasible for larger values of l .

Instead consider the set of all integers which are within a Hamming distance d from S_1 . The total number of such integers is $\binom{l}{d} \times 2^d$. Note that the solution set is a subset of this set. Thus we can search through this set for the solution, significantly reducing our search space. This search space can be enumerated starting with bit representation of S_1 and inverting bits in selected positions. Specifically, we traverse this l -bit vector from left to right. Suppose we are at position i ($0 \leq i < l - 1$) and we have inverted j positions so far. If $j < d$, we have two possibilities for position $i + 1$, either to invert or keep the bit the same. We consider both possibilities and move forward. If $j = d$, we don't consider the possibility of inverting the bit. Once $i = l - 1$, we check the Hamming distance of the resulting number from each element in the set $\{S_2 \dots S_n\}$. If all the distances are less than or equal to d , the number goes into the solution set. One can come up with more sophisticated strategies to further reduce the search space and you are welcome (but not required) to do so.