**Automated System for Structural Break Detection in DNA Sequences**

# Manual (v1.0)

## Background

ASSBD stands for Automatic System for Structural Break Detection. This is an automated system that will identify the DNA sequences breaks. Breaks includes space, comma and any other symbols. This also helps to determine the mutation, mismatches, length of the common segments. Consecutively, we can find out the any repeated segment by providing length of the data set. It will help to find out any desired length repeated segment.

## Availability and implementation

Download latest version of ASSBD.jar from

https://github.com/abmmki/assbd

Released under GNU public license version 3 (GPL v3).

**Contact:** khademul@du.ac.bd**,** sarwar.saubdcoxbazar@gmail.com

## Installation

The program needs to have Java Development Kit (J.D.K) and JRE 1.7 (or latest version) preinstalled. OS requirement is Windows 7. This program also tested in Linux (Ubuntu) OS.

The basic computer requirement is Pentium 4 and any update versions will do. Basic primary requirement is 1GB and secondary memory (Hard Disk) is 40BG.

# Chapter one: Sequence Alignment

Sequence Alignment can be performing in both local and global concept. At first we have to double click on the system named ASSBD to open **(Figure 1).**
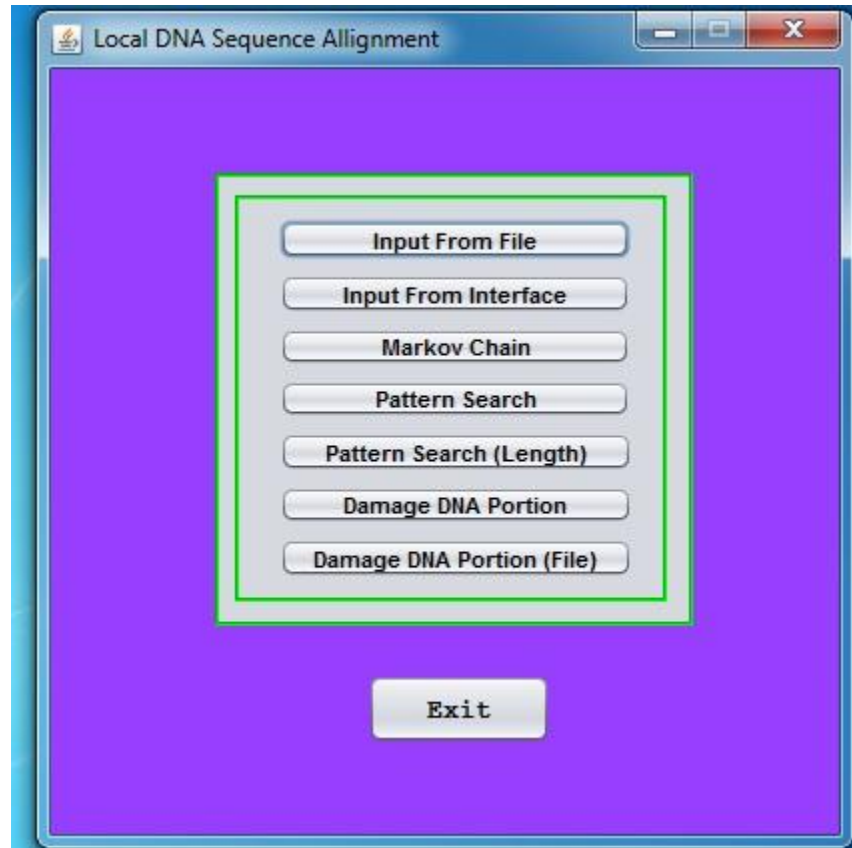


**Figure 1**: ASSBD Complete Interface

The process is give below in step by step process

1. Click on Input from File and new interface (Figure 2) will visible. Here we see **File name** that is data file. The **Browse** button enables us to select the file from certain position of computer. Two position buttons exist as **First position** that is from where we want to start the alignment and **Last position** indicates the value of the sequences up to length we want to align.

**Figure 2**: Input from file interface

2. Click on Browse button. New interface (**Figure 3**) will open that will enable data set from any part of the computer.



**Figure 3**: Browse Interface for data file selection

3. An example using the data set from desktop has performed (**Figure 4 and Figure 5**).



**Figure 4**: A data set file has selected from desktop, named Data.txt



**Figure 5:** Result after data set selection.

4. Same operation as 1 and 2 can be done by using keyboard input by clicking at input from keyboard button (**Figure 1**). In response of that new window (Figure 6) will open.



**Figure 6:** Input data set using Keyboard

5. An example using keyboard data set (**Figure 7**) is narrated.

**Figure 7**: An impact of Keyboard

6. Consequently, the pattern search is also an important part of ASSBD. When we click on the pattern search button (**Figure 8**) and we get the outcome (**Figure 9**). This enables the users to find specific pattern from given data set. That is, which pattern is repeatedly exist, will show the outcome. For example, given data set (**Figure 9**) has illustrated that **AGA** pattern is 21% of the total data set and its frequency is 6, i.e. there are six times **AGA** pattern.

Figure 8:  Outline for Pattern Search

Figure 9: The resultant on pattern search

7. Here, we also find the repeated pattern by providing the length of the pattern (**Figure 10**). The outcome for a certain data set, we have noticed a four length pattern **AAGG** is resulted (**Figure 11**).

**Figure 10**: Pattern search according to the length



Figure 11: The certain patter finding based on length.

# Chapter 2: Structural Break

Structural Break is the changes in DNA sequences as any kind of symbols or space among sequences.

AAAA-GCT C-CCCTTTTAT-A-CTGATGCTG-A
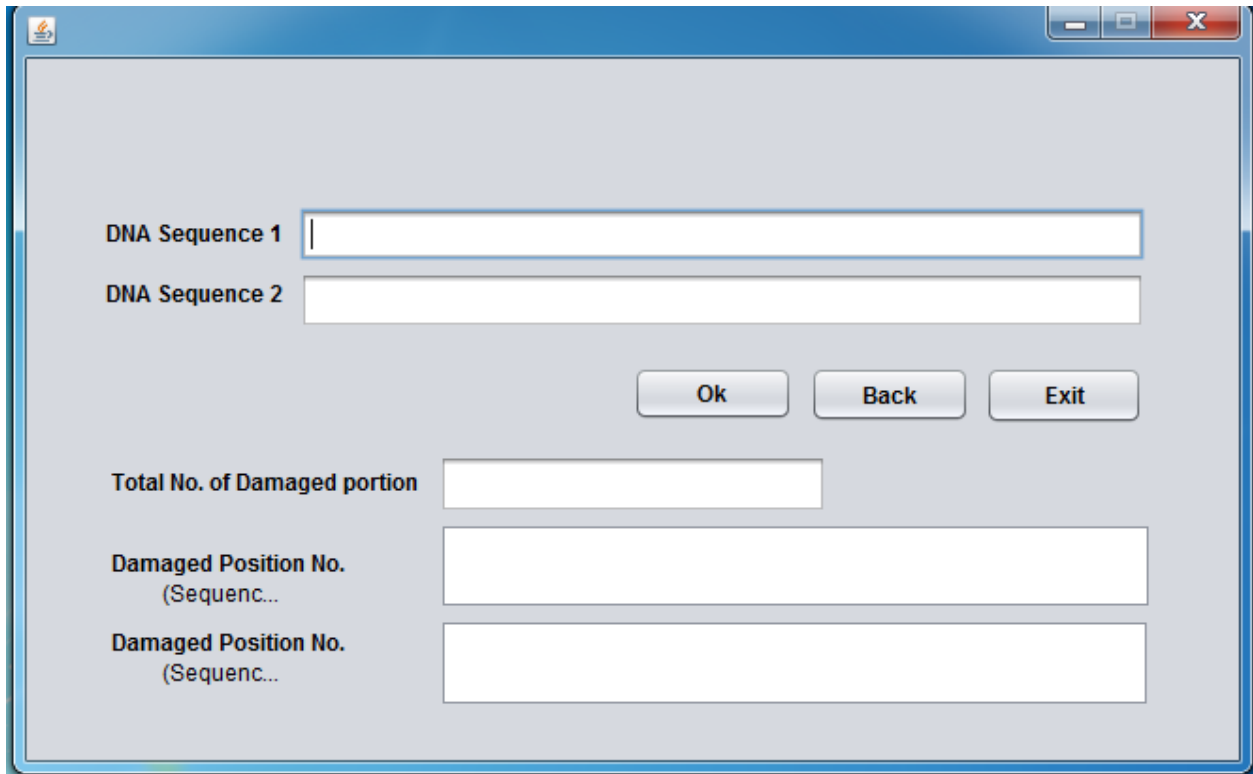
Here, there are structural break at four points.



Figure 12: Opening face for Structural Break

Figure 13: Outcome for Given Data set

It is possible to get input by making browsing data set from computer. This development (Figure 14) helps easy browsing for analysis.

Figure 14: Interface for file input for damage selection and outcome shows (Figure 15)

Figure 15: Outcome for file browsing data set. Here we see that total 15 structural breaks occurred and the positions are identified in both sequences

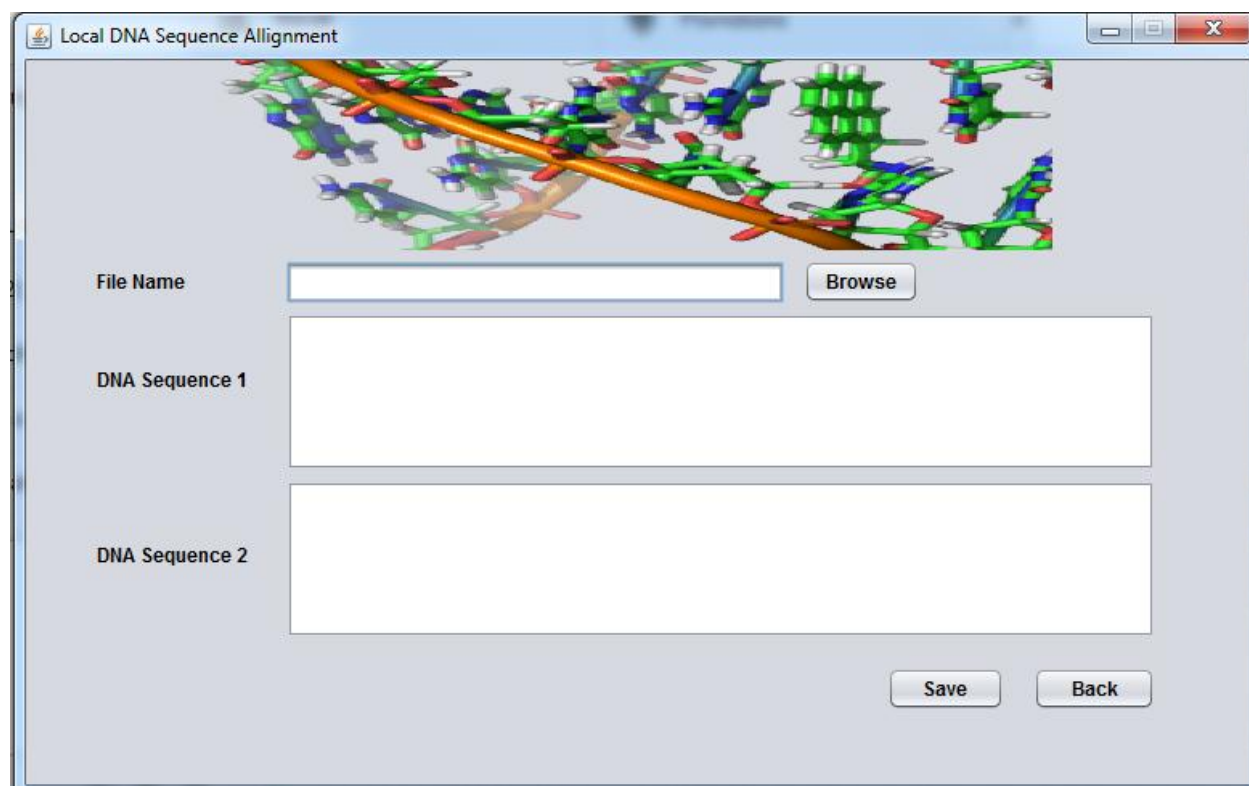It is also possible to change the data set from the interface (Figure 16). For this user have to click change file (Figure 15)

Figure 15: Data set change options.