Project 4: Gene Sequencing

1. Time and complexity of methods

```
# def unrestricted prep array
# Preps array by adding first row
# Time Complexity: O(n) Creates n cells
# Space Complexity: O(n) Stores n cells
# def unrestricted algorithm
# Sequences 2 strings
# Time Complexity: O(n * m) For each letter in string 1, traverse string 2
# Space Complexity: O(n * m) Makes and stores things in array of size n * m
# def banded_prep_array
# Preps array by adding first row
# Time Complexity: O(k) Creates k cells, k = 4
# Space Complexity: O(k) Stores k cells in list, k = 4
# def banded algorithm
# Sequences 2 strings
# Time Complexity: O(n * k) For each letter in string 1, visit 4 times, k = 4
# Space Complexity: O(n * k) Makes and stores things in array of size n with 4 cells, k = 4
# def get smallest combination
# returns smallest combination
# Time Complexity: O(n) visit n nodes and make a string
# Space Complexity: O(n) saves a string at about n long
# def get score
# returns the score
# Time Complexity: O(1) Just returns the score
# Space Complexity: Not really applicable
```

2. Paragraph

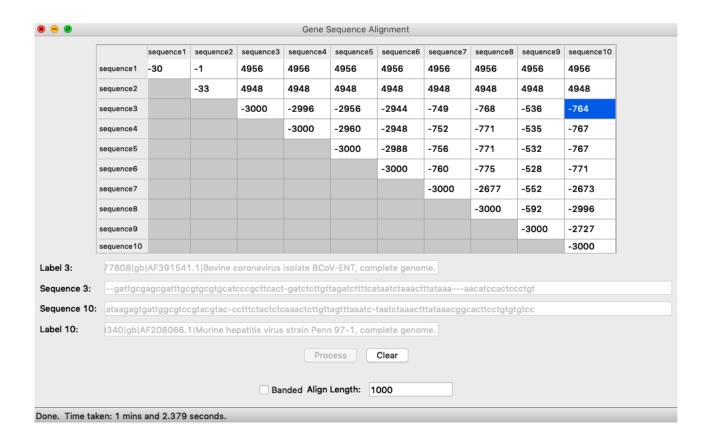
```
# KEY Array [i][j]
#
#
    String 1
#
# S jjjjjjjjj
# t i
# r i
# i i
# n i
#gi
# 2 i
#
\# diagonal_index = [i - 1] [j - 1]
# up index
            = [i - 1] [j]
           = [i] [j-1]
# left index
```

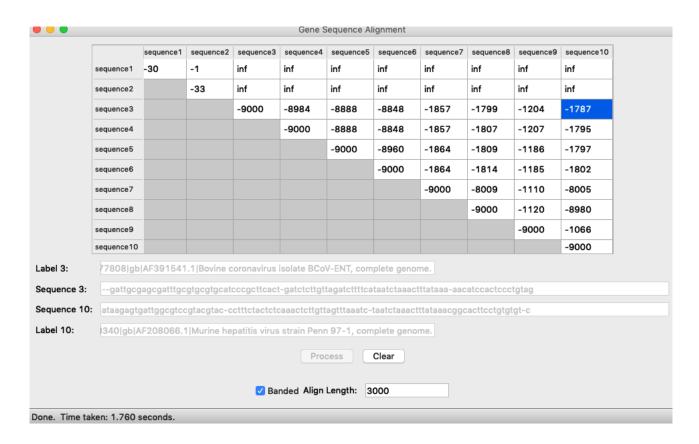
My algorithm start at index [0][0] and makes the first row of insertions, deletions to the right. Then for each subsequent row I reference the left child, upper diagonal child and upper child. I calculate the costs of each parent, determine the transition with the least amount of cost, then store that node and the new cells parent.

Because I am only storing the upper half diagonal of the data to avoid duplication, I did self.string1_len + 1 - (i - 1) loops on each subsequent row. So basically n - 1 loops for each additional row.

Because I was only storing half the data and in a list format, tracking the parents was a little difficult. To access a diagonal parent, I accessed the cell right above. A upper parent was the cell right above and one the left. The left parent was the last cell in the list I was adding to.

3. Results





```
# KEY Array [i][i]
#
#
   String 1
#
# S jjjjjjjjj
# t i
#ri
# i i
# n i
#gi
# 2 i
#
from enum import Enum
# Cell
# Stores data for each cell in array
class Cell:
 def init (self, parent, transition type, cost):
   self.parent = parent
   self.transition_type = transition_type
   self.cost = cost
# Transition
# Stores different transitions
class Transition(Enum):
 UP = 0
 LEFT = 1
 DIAGONAL = 2
# Sequence
# Sequences 2 strings
# Time Complexity: Depends on user input
# if banded see banded algorithm complexity
# if not see unrestricted algorithm complexity
# Space Complexity: Depends on user input
# if banded see banded_algorithm complexity
# if not see unrestricted algorithm complexity
class Sequence:
 def __init__(self, string1, string2, banded):
   self.string1 = string1
   self.string1_len = len(string1)
   self.string2 = string2
   self.string2_len = len(string2)
   self.array = []
```

```
if not banded:
    self.unrestricted_prep_array()
    self.unrestricted algorithm()
    self.score = self.get score()
  else:
    # We do not allow more than 3 insertions/ deletion
    if self.string1_len - self.string2_len <= 4:</pre>
      self.banded_prep_array()
      self.banded_algorithm()
      self.score = self.get score()
    else:
      self.score = float('inf')
# def unrestricted prep array
# Preps array by adding first row
# Time Complexity: O(n) Creates n cells
# Space Complexity: O(n) Stores n cells
def unrestricted prep array(self):
  # add first row to array
  self.array.append([])
  # temp cell
  tmp cell = Cell(None, None, 0)
  # add first cell to array
  self.array[0].append(tmp_cell)
  # add remaining cells to row
  for j in range(1, self.string1 len + 1):
    tmp cell = Cell(tmp cell, Transition.LEFT, j * 5)
    self.array[0].append(tmp cell)
# def unrestricted algorithm
# Sequences 2 strings
# Time Complexity: O(n * m) For each letter in string 1, traverse string 2
# Space Complexity: O(n * m) Makes and stores things in array of size n * m
def unrestricted algorithm(self):
  # for each row (string 2)
  for i in range(1, self.string2_len + 1):
    # add another row to array
    self.array.append([])
    for j in range(1, self.string1 len + 1 - (i - 1)):
      \# diagonal_index = [i - 1] [j - 1]
      # up index = [i - 1] [j]
      # left index
                     = [i] [j-1]
      # get up cost
      up_cost = self.array[i - 1][i].cost + 5
      # get left cost
```

```
if len(self.array[i]) > 0:
         left_cost = self.array[i][-1].cost + 5
       else:
         left cost = float("inf")
       if self.string1[j - 1 + (i - 1)] == self.string2[i - 1]:
         strings match = True
       else:
         strings_match = False
       if strings match:
         # strings are the same, cost: -3
         diagonal_cost = self.array[i - 1][j - 1].cost - 3
         # strings are not the same, cost: +1
         diagonal_cost = self.array[i - 1][j - 1].cost + 1
       # stores smallest parent
       if diagonal_cost <= up_cost and diagonal_cost <= left_cost:</pre>
         if strings_match:
           self.array[i].append(Cell(self.array[i - 1][i - 1],
                           Transition.DIAGONAL.
                           diagonal cost))
         else:
           self.array[i].append(Cell(self.array[i - 1][j - 1],
                           Transition.DIAGONAL.
                           diagonal_cost))
       elif up_cost <= left_cost:</pre>
         self.array[i].append(Cell(self.array[i - 1][j],
                         Transition.UP,
                        up cost))
       else:
         self.array[i].append(Cell(self.array[i][-1],
                         Transition.LEFT,
                        left cost))
# def banded_prep_array
# Preps array by adding first row
# Time Complexity: O(k) Creates k cells, k = 4
# Space Complexity: O(k) Stores k cells in list, k = 4
def banded prep array(self):
  # add first row to array
  self.array.append([])
  # temp cell
  tmp_cell = Cell(None, None, 0)
  # add first cell to array
  self.array[0].append(tmp_cell)
  # add remaining cells to row
  for i in range(1, 4):
    tmp_cell = Cell(tmp_cell, Transition.LEFT, j * 5)
    self.array[0].append(tmp_cell)
# def banded_algorithm
```

```
# Sequences 2 strings
# Time Complexity: O(n * k) For each letter in string 1, visit 4 times, k = 4
# Space Complexity: O(n * k) Makes and stores things in array of size n with 4 cells, k = 4
def banded_algorithm(self):
  # for each row (string 2)
  for i in range(1, self.string2_len + 1):
     # add another row to array
     self.array.append([])
     # for each column (string 1)
     for j in range(1, 5):
       if self.string1_len - i < 3 and j == 4:
       if self.string1_len - i < 2 and j == 3:
          continue
       if self.string1 len - i < 1 and i == 2:
          continue
       if self.string1_len - i < 0 and j == 1:
          continue
       \# diagonal_index = [i - 1] [j - 1]
                        = [i - 1] [j]
       # up index
       # left_index
                        = [i] [j-1]
       # get up cost
       if i < 4:
          up_cost = self.array[i - 1][j].cost + 5
       else:
          up_cost = float("inf")
       # get left cost
       if len(self.array[i]) > 0:
          left_cost = self.array[i][-1].cost + 5
       else:
          left_cost = float("inf")
       if self.string1[j - 1 + (i - 1)] == self.string2[i - 1]:
          strings match = True
       else:
          strings_match = False
       # get diagonal cost
       if strings match:
          # strings are the same, cost: -3
          diagonal_cost = self.array[i - 1][j - 1].cost - 3
          # strings are not the same, cost: +1
          diagonal_cost = self.array[i - 1][j - 1].cost + 1
       # stores smallest parent
       if diagonal_cost <= up_cost and diagonal_cost <= left_cost:</pre>
          if strings_match:
             self.array[i].append(Cell(self.array[i - 1][j - 1],
                             Transition.DIAGONAL,
                             diagonal_cost))
```

```
self.array[i].append(Cell(self.array[i - 1][j - 1],
                          Transition.DIAGONAL,
                          diagonal cost))
      elif up_cost <= left_cost:</pre>
         self.array[i].append(Cell(self.array[i - 1][j],
                        Transition.UP,
                       up cost))
      else:
         self.array[i].append(Cell(self.array[i][-1],
                        Transition.LEFT.
                       left cost))
# def get smallest combination
# returns smallest combination
# Time Complexity: O(n) visit n nodes and make a string
# Space Complexity: O(n) saves a string at about n long
def get combo string(self):
  if self.score == float('inf'):
    return ["No Alignment Possible", "No Alignment Possible"]
  combo string a = ""
  combo_string_b = ""
  tmp cell = self.array[self.string2 len][-1]
  string1_index = self.string1_len - 1
  string2 index = self.string2 len - 1
  while tmp cell is not None:
    if tmp_cell.transition_type == Transition.UP:
      combo_string_a = "-" + combo_string_a
      combo_string_b = self.string2[string2_index] + combo_string_b
      string2_index -= 1
    elif tmp_cell.transition_type == Transition.LEFT:
      combo string a = self.string1[string1 index] + combo string a
      combo string b = "-" + combo string b
      string1 index -= 1
    elif tmp_cell.transition_type == Transition.DIAGONAL:
      combo_string_a = self.string1[string1_index] + combo_string_a
      combo_string_b = self.string2[string2_index] + combo_string_b
      string1 index -= 1
      string2_index -= 1
    tmp cell = tmp cell.parent
  return [combo_string_b[:100], combo_string_a[:100]]
# def get_score
# returns the score
# Time Complexity: O(1) Just returns the score
# Space Complexity: Not really applicable
```

else:

def get_score(self):
 return self.array[self.string2_len][-1].cost