

Senior Project Documentation Title

First Student

Bachelor of Science in Information Technology

Second Student

Bachelor of Science in Computer Science

Third Student

Bachelor of Science in Information Systems

Senior project submitted to the faculty of the
Department of Computer Science
College of Computer Studies, Ateneo de Naga University
in partial fulfillment of the requirements for their respective
Bachelor of Science degrees

Project Advisor: Adviser's Name

First Panel Member

Second Panel Member

Third Panel Member

Month Day, 2020

Naga City, Philippines

Keywords: keywordOne, keywordTwo, keywordThree

Copyright 2020, First Student, Second Student, and Third Student

The Senior Project entitled

Senior Project Documentation Title

developed by

First Student

Bachelor of Science in Information Technology

Second Student

Bachelor of Science in Computer Science

Third Student

Bachelor of Science in Information Systems

and submitted in partial fulfillment of the requirements of their respective Bachelor of Science degrees
has been rigorously examined and recommended for approval and acceptance.

First Panel Member

Panel Member

Date signed: _____

Second Panel Member

Panel Member

Date signed: _____

Third Panel Member

Panel Member

Date signed: _____

Adviser's Name

Project Advisor

Date signed: _____

The Senior Project entitled

Senior Project Documentation Title

developed by

First Student

Bachelor of Science in Information Technology

Second Student

Bachelor of Science in Computer Science

Third Student

Bachelor of Science in Information Systems

and submitted in partial fulfillment of the requirements of their respective Bachelor of Science degrees is hereby approved and accepted by the Department of Computer Science, College of Computer Studies, Ateneo de Naga University.

Marianne P. Ang, MS

Chair, Department of Computer Science

Date signed: _____

Joshua C. Martinez, MIT

Dean, College of Computer Studies

Date signed: _____

Declaration of Original Work

We declare that the Senior Project entitled

Senior Project Documentation Title

which we submitted to the faculty of the

Department of Computer Science, Ateneo de Naga University

is our own work. To the best of our knowledge, it does not contain materials published or written by another person, except where due citation and acknowledgement is made in our senior project documentation. The contributions of other people whom we worked with to complete this senior project are explicitly cited and acknowledged in our senior project documentation.

We also declare that the intellectual content of this senior project is the product of our own work. We conceptualized, designed, encoded, and debugged the source code of the core programs in our senior project. The source code of third party APIs and library functions used in my program are explicitly cited and acknowledged in our senior project documentation. Also duly acknowledged are the assistance of others in minor details of editing and reproduction of the documentation.

In our honor, we declare that we did not pass off as our own the work done by another person. We are the only persons who encoded the source code of our software. We understand that we may get a failing mark if the source code of our program is in fact the work of another person.

First Student

4 - Bachelor of Science in Information Technology

Second Student

4 - Bachelor of Science in Computer Science

Third Student

4 - Bachelor of Science in Information Systems

This declaration is witnessed by:

Adviser's Name

Project Advisor

Senior Project Documentation Title

by

First Student, Second Student, and Third Student

Project Advisor: Adviser's Name

Department of Computer Science

EXECUTIVE SUMMARY

The Shortest Common Superstring (SCS) problem, known to be NP-Complete, seeks the shortest string that contains all strings from a given set. This paper provides a detailed summary of SCS problem including the problem's complexity and approximability, and shows that DNA sequence assembly can be approached using SCS. Then this paper discusses the two most popular approximation algorithms for the SCS problem namely greedy and genetic algorithm. Finally, this study presents results supporting the simplicity and performance guarantee of the greedy algorithm as well as the convergence of the genetic algorithm to satisfactory approximate solutions as individuals representing problem instances evolve.

I dedicate this research work to all of humanity.

ACKNOWLEDGEMENTS

I thank everyone who helped me finish this thesis.

TABLE OF CONTENTS

1	Introduction	1
1.1	Project Context	2
1.2	Purpose and Description	2
1.3	Objectives	2
1.4	Scope and Limitations	3
2	Review of Related Systems and Related Literature	4
2.0.1	Greedy-SCS	5
2.0.2	Genetic Algorithm	5
3	Technical Background	7
4	Methodology	9
4.1	Systems Analysis	9
4.2	Systems Design	9
4.3	Requirements Specification	9
4.4	Development and Testing	9
5	Results	10
6	Contributions and Recommendations	13
6.1	Summary of Contributions	13
6.2	Recommendations	13
A	Code Listing	15

LIST OF FIGURES

1.1	Sample overlap graph with each adjacent nodes having at least $k = 3$ overlaps. The original string is GCATTATATATTGCGCGTACGGCGCCGCTACA.	2
3.1	SCS Decision Problem	8
5.1	Best superstring as a function of generations. Each point in the figure represents the average of 50 runs on 50 different randomly generated problem instances. For each such instance, two runs were performed, the better of which was considered for statistical purposes	11

LIST OF TABLES

5.1	Average superstring length and runtime of our Genetic Algorithm on 50 randomly generated problem instances. Each input set contains 50 blocks.	11
-----	--	----

Chapter 1

Introduction

The Shortest Common Superstring (SCS) problem, known to be NP-Complete, seeks the shortest string that contains all strings from a given set. In this paper, we provide the summary of the problem and some of its characteristics.

The SCS problem has been extensively studied for its applications in string compression and DNA sequence assembly [7].

The superstring problem has applications to data storage, specifically, data (string) compression [5]. In many programming languages, a character string may be represented by a pointer to that string. The problem for the compiler is to arrange strings so that they may be “overlapped” as much as possible.

DNA sequence assembly is another problem to which an SCS algorithm is known to apply. The *sequencing* problem in molecular biology is to “read” a string of DNA, which can be viewed as a string over the alphabet $\{A,C,G,T\}$. Sequencing produces such a large number of fragments that almost all genome positions are covered by many fragments. This short fragments thus have large overlaps between other pieces. Hence, they can be given as an input to SCS algorithm. Figure 1.1 shows an overlap graph consisting DNA reads (or fragments) as nodes.

In [7]’s paper, SCS was used to analyze DNA sequence assembly using a greedy algorithm.

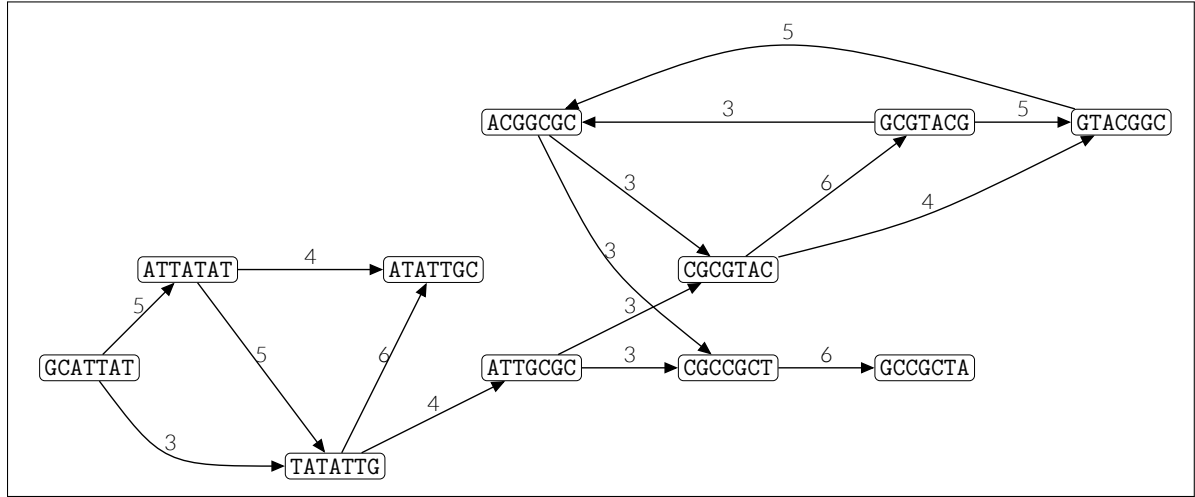


Figure 1.1: Sample overlap graph with each adjacent nodes having at least $k = 3$ overlaps. The original string is GCATTATATATTGCGGTACGGCGCCGCTACA.

1.1 Project Context

1.2 Purpose and Description

1.3 Objectives

1.4 Scope and Limitations

Chapter 2

Review of Related Systems and Related Literature

Several researches have provided in depth study on approximation algorithms for the SCS problem.

Most [11, 7, 12], however, considered a greedy algorithm to solve the SCS problem. The reason why greed works for shortest common superstring problem was explored by [7]. They explained the good performance of the greedy algorithm by using the smoothed analysis. That is, for any given instance I of *SCS*, the average approximation ratio of the greedy algorithm on a small random perturbation of I is $1 + o(1)$.

Apart from the greedy approach, other methods including bio-inspired algorithms, have also been employed to solve lots of hard problems in computer science including the SCS problem. In [12]’s work, for example, four approaches for finding solutions to the SCS problem: a standard genetic algorithm, a novel cooperative-coevolutionary algorithm, a benchmark greedy algorithm, and a parallel coevolutionary-greedy approach have been compared. In the paper, the coevolutionary approach produced the best results.

Zaritsky et al’s [12] work served as the foundation of this paper. Hence, the following algorithms, which have been explored and tested in the latter, have been extracted and validated whether they are indeed generating good approximate solution to the SCS problem.

2.0.1 Greedy-SCS

The core idea of the greedy algorithm is to repeatedly merge pairs of distinct strings with maximum overlap until only one remains. It has been conjectured by [8] that the superstring produced by the greedy algorithm is at most two times the optimal. Below is a pseudocode for the greedy algorithm:

Algorithm 1 GREEDY(R)

Input: R - set of strings

Output: superstring of set R

- 1: **while** $|R| > 1$ **do**
 - 2: choose $x_1 \neq x_2 \in R$ such that $overlap(x_1, x_2)$ is maximal
 - 3: $R \leftarrow (R - \{x_1, x_2\}) \cup \{merge(x_1, x_2)\}$
 - 4: **end while**
 - 5: **return** remaining string in R
-

The greedy algorithm is probably the most widely used heuristic used in DNA-sequencing due to its simplicity and constant performance guarantee which is most likely to be factor-2.

2.0.2 Genetic Algorithm

In this paper, we follow the standard model for the genetic algorithm (see Algorithm 2). That is, the algorithm generates an initial population of random candidate solutions, then the process of selection based on a defined fitness function, crossover, and mutation are used to evolve the next generation, of which each individual is then evaluated and assigned a fitness value. These steps are repeated a predefined number of times or until the solution is satisfactory.

Algorithm 2 GENERIC GA()

```
1:  $t \leftarrow 0$  {generation number}
2: initialize Population $t$ 
3: evaluate(Population $t$ )
4: while termination condition not met do
5:   select individuals from Population $t$ 
6:   recombine individuals
7:   mutate individuals
8:   Population $t+1$   $\leftarrow$  newly created individuals
9:    $t \leftarrow t + 1$ 
10:  evaluate(Population $t$ )
11: end while
12: return solution derived from the best individual in Population $t$ 
```

Chapter 3

Technical Background

A superstring is simply a string over some alphabet for which given a set of string from the same alphabet, the latter's members are all substrings of the former. To understand the SCS problem, it is best to assume some important concepts related to the problem itself. First, let us assume some set $R = \{x_1, \dots, x_n\}$ as a set of strings (or blocks) over some alphabet Σ . Formally, we define a *superstring* of R as a string w containing each $x_i \in R$, as a substring.

Following are the elementary operations associated with the construction of a superstring. The *overlap*(u, v) of two strings u and v is the maximum overlap between u and v . That is to say, the longest suffix of u (in terms of characters) that is a prefix of v . The *prefix*(u, v) of u and v is the prefix of u obtained by removing its overlap with v . Lastly, *merge*(u, v) is the concatenation of u and v with the overlap appearing only once.

As an example, say we have, $\Sigma = \{a, b, c\}$ and $R = \{cbcaca, cacac\}$. The string *cacacbcaca* is a superstring of R while the string *cbcacac* is the shortest common superstring. The following relations also hold:

$$\begin{aligned} \text{overlap}(cbcaca, cacac) &= caca, \\ \text{overlap}(cacac, cbcaca) &= c, \\ \text{prefix}(cbcaca, cacac) &= cb, \\ \text{merge}(cbcaca, cacac) &= cbcacac. \end{aligned}$$

A *superstring* w is also defined as the string $\text{prefix}(x_1, x_2) \cdot \text{prefix}(x_2, x_3) \cdots \text{prefix}(x_n, x_1) \cdot \text{overlap}(x_n, x_1)$. That is, a *superstring* is simply a concatenation of all the strings “minus” the overlapping duplicates. Apparently, each superstring of a set of strings defines a permutation of

the set's elements. Conversely, every permutation of the set's elements corresponds to a single superstring

Our interest however lies on defining the SCS problem. Essentially, the SCS problem is: Given a set of strings R and a positive integer K , does R have a superstring of length K ?

Figure 3.1 shows the SCS Decision Problem following the template provided by Garey and Johnson [6].

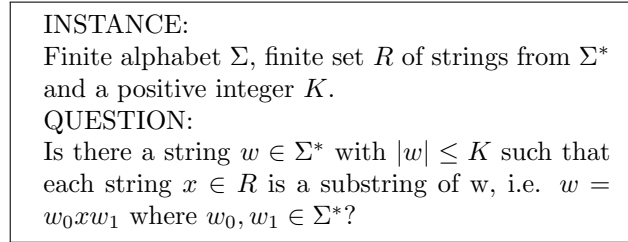


Figure 3.1: SCS Decision Problem

In the Compendium of NP Optimization problems published by [3] and in the list of NP-Complete Problems published by [6] the SCS problem appears under under Storage and Retrieval (SR).

Chapter 4

Methodology

Methodology stuff will appear here.

4.1 Systems Analysis

4.2 Systems Design

4.3 Requirements Specification

4.4 Development and Testing

Chapter 5

Results

The researcher conducted an experiment to validate the results in [12]’s work. Following is the description of the experiment done.

The input strings used in the experiment were generated in a manner similar to the one used in DNA sequencing. That is, a random string is generated, duplicated a predetermined number of times, and the copies are randomly divided into blocks of a given size. The set of all these blocks is the input to the SCS problem. The reasons for choosing this type of input has been extensively explained in [12]. Since our random string is a sequence of characters A,C,G, and T, we must transform this into an equivalent binary string. To do this, we simply assign each letter a unique two-bit representation, that is, A = 00, C = 01, G = 11, and T = 10.

The parameters used in the input generation is as follows:

Size of random string: 250 bases or characters

Minimal block size: 20 characters

Maximal block size: 30 characters

Number of duplicates created from the random string: 5

Let l denote the length of the derived string in the genetic algorithm case. Let m denote the number of blocks not covered by the derived string, and let b denote the maximal block size (30, in our case). The fitness value, f , of an individual is computed as follows:

$$f = \frac{1}{l + m * b}$$

This fitness function above drives evolution towards shorter superstrings covering as many blocks as

Number of Generations	Average Superstring Length (Rounded to Nearest Units)	Average Runtime (in Seconds)
100	760	34.57
250	596	74.20
500	526	139.70
1000	482	267.09
5000	421	1038.73
10000	413	1943.55
20000	409	3700.57

Table 5.1: Average superstring length and runtime of our Genetic Algorithm on 50 randomly generated problem instances. Each input set contains 50 blocks.

possible.

We performed 50 randomly generated problem instances. On each problem instance each type of genetic algorithm was executed twice and the better run of the two was used for statistical purpose. The results are summarized in Table 5.1 and is presented graphically in Figure 5.1.

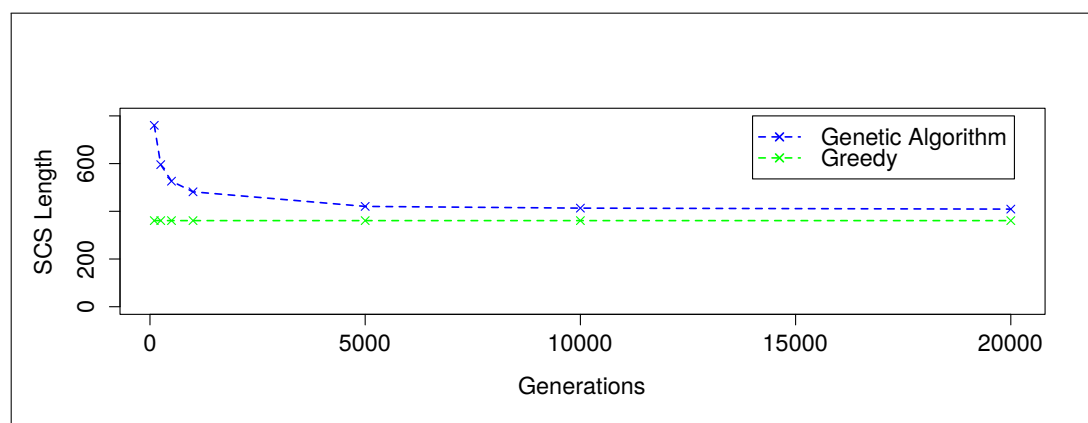


Figure 5.1: Best superstring as a function of generations. Each point in the figure represents the average of 50 runs on 50 different randomly generated problem instances. For each such instance, two runs were performed, the better of which was considered for statistical purposes

Note that our results seem to disagree with that of [12]. In contrast with the latter's findings, our tests shows that the greedy algorithm outperformed the genetic algorithm both in terms of finding superstrings with minimal length and the the time it takes to do so.

On the same set of test cases, the greedy algorithm produced superstrings having an average of 361 characters, in an average of runtime 130.64 milliseconds. This is far better than the 409 average

superstring length that the genetic algorithm derived in 20000 generations in more or less an hour of operation. Note that the average length of the superstring derived by the greedy algorithm is less than the conjectured factor-2 performance guarantee. Recall that the length of the initial (reference) string is 250.

Still it is good to note that the genetic algorithm approach in solving the SCS problem converges. That is, we see a better solution yield as time or generations progress.

Chapter 6

Contributions and Recommendations

6.1 Summary of Contributions

6.2 Recommendations

There could be several reasons why our tests did not yield the “same” results as that of [12]’s. The author claims that this may be due to the myriad of parameters used in evolving individuals in each step of the genetic algorithm and the large number of possible candidates or values that we can set these parameters with. For example, we may claim that the appropriate choice of the recombination technique is necessary. In [12], a two-point crossover was employed. Thinking that this may be less of a factor, the author proceeded with just one-point cross in implementing the genetic algorithm for the experiment. Another could be in choice of the probabilities associated with recombination, elitism, and mutation rates. The author did only consider what was stated in the base paper without performing tests to validate whether this values are appropriate and could still complement the change made in some of the assumptions. As such, more empirical tests could be done to appropriately estimate these parameters.

If indeed DNA sequence assembly is the task at hand, lots of other techniques are now available. One may consider for example representing the reads or DNA fragments in a de Bruijn graph as

opposed to the overlap graph presented in this paper. Assembling reads in a de Bruijn graph reduces the problem to a fragment assembly problem that can be formulated as the goal to find a trail or Eulerian path that visits each edge (read or fragment) in the (de Bruijn) graph exactly once. Such, somehow makes the assembly process much “easier” since Eulerian path construction is known to be solvable in deterministic polynomial time.

Appendix A

Code Listing

```
1  /**
2   * File:    scs-greedy.cpp
3   * Author:  glenn
4   *
5   * Created on May 13, 2015, 10:52 AM
6   */
7
8  #include <cstdlib>
9  #include <iostream>
10 #include <fstream>
11 #include <string>
12 #include <list>
13
14 using namespace std;
15
16 /**
17  * Greedy SCS
18  */
19 string find_max_overlap(string source, string sink);
20 void showR(list<string> R){
21     for(list<string>::iterator r = R.begin(); r != R.end(); r++) {
22         cout << *r << endl;
23     }
24 }
25 void greedy_scs(list<string> &R);
26
27 // remove substrings
28 void reduce(list<string> &R){
29     list<string> S;
30     for (list<string>::iterator x = R.begin(); x != R.end(); x++) {
31         bool is_substr = false;
32         for (list<string>::iterator y = R.begin(); y != R.end(); y++) {
33             if(*x == *y ) continue;
34             if((*y).find(*x) != -1){
35                 cout << *x << " " << *y << endl;
36                 is_substr = true;
37             }
38         }
39     }
40 }
```

```

38     }
39     if(!is_substr){
40         S.push_back(*x);
41     }
42 }
43 R.clear();
44 R.insert(R.end(), S.begin(), S.end());
45 }
46
47 int main(int argc, char** argv) {
48
49     if(argc != 2){
50         cout << "Usage: " << argv[0] << " <input-file>" << endl;
51         return -1;
52     }
53
54     int accessions = 50;
55     {
56         list<string> R;
57         string reference;
58         string fragment;
59
60         string filename = argv[1];
61         ifstream fin(filename.c_str());
62         if(!fin){
63             cout << "Could not open file \' " << filename << "\' " << endl;
64             return -1;
65         }
66
67         fin >> reference;
68         while(fin >> fragment){
69             R.push_back(fragment);
70         }
71
72         //~ R.push_back("AACTG");
73         //~ R.push_back("AACA");
74         //~ R.push_back("ATC");
75         //~ R.push_back("ATCCG");
76         //~ R.push_back("AAT");
77         //~ R.push_back("AATA");
78         //~ R.push_back("AATC");
79         //~ R.push_back("AGG");
80
81         //~ reduce(R);
82
83         int start_s=clock();
84         //~ the code you wish to time goes here
85         greedy_scs(R);
86         int stop_s=clock();
87         //~ cout << "Time: " << (stop_s-start_s)/double(CLOCKS_PER_SEC)*1000 << endl;
88
89         //~ cout << reference << endl;
90         //~ cout << reference.length() << endl;
91         //~ cout << endl;
92
93         cout << R.front().length() << ", "
94             << (stop_s-start_s)/double(CLOCKS_PER_SEC)*1000 << ", "
95             << R.front() << endl;

```

```

96  }
97
98      return 0;
99  }
100
101  string find_max_overlap(string source, string sink) {
102
103      int k = min(source.length(), sink.length()); // search windows
104      string suffix = source.substr(source.length() -
105          k, source.length()); // last k characters
106      string prefix = sink.substr(0, k); // first k characters
107      string matched = "";
108      while (true) {
109          if (suffix == prefix)
110              return suffix;
111          else {
112              k--;
113              suffix = source.substr(source.length() - k, source.length());
114              prefix = sink.substr(0, k);
115          }
116      }
117      return "";
118  }
119
120  void greedy_scs(list<string> &R) {
121      // find the maximum overlap strings
122
123      string overlap_xy, max_overlap_x, max_overlap_y;
124      int max_overlap = 0;
125      int iter = 0;
126      while (R.size() > 1) {
127
128          iter++;
129          max_overlap = 0;
130          overlap_xy = max_overlap_x = max_overlap_y = "";
131          //~ cout << "finding overlap: R = " << R.size() << " t = " << iter << endl;
132          //~ showR(R);
133
134          for (list<string>::iterator x = R.begin(); x != R.end(); x++) {
135              for (list<string>::iterator y = R.begin(); y != R.end(); y++) {
136
137                  if (x == y)
138                      continue;
139
140                  overlap_xy = find_max_overlap(*x, *y);
141
142                  if (overlap_xy.length() >= max_overlap) {
143                      max_overlap = overlap_xy.length();
144                      max_overlap_x = *(x);
145                      max_overlap_y = *(y);
146                  }
147              }
148          }
149      }
150
151      // cout << "max overlap: " << max_overlap << ": "
152      // << max_overlap_x << ", " << max_overlap_y << endl;

```

```
155
156         // remove the two strings
157
158         for (list<string>::iterator it = R.begin(); it != R.end(); it++) {
159             if (*it == max_overlap_x) {
160                 it = R.erase(it);
161             }
162             if (*it == max_overlap_y) {
163                 it = R.erase(it);
164             }
165         }
166
167         // merge
168         string merged_xy;
169         //~ cout << "maxoverlap : " << max_overlap << endl;
170
171         //~ cout << endl;
172         //~if(max_overlap != 0)
173             merged_xy = max_overlap_x.substr(0, max_overlap_x.length() - max_overlap) + m
174         //~ else
175             //~ merged_xy = max_overlap_x + max_overlap_y;
176
177         if (merged_xy != "")
178             R.push_back(merged_xy);
179
180     }
181
182
183 }
```

REFERENCES

- [1] C. ARMEN AND C. STEIN, *Improved length bounds for the shortest superstring problem*, in Algorithms and Data Structures, S. Akl, F. Dehne, J.-R. Sack, and N. Santoro, eds., vol. 955 of Lecture Notes in Computer Science, Springer Berlin Heidelberg, 1995, pp. 494–505.
- [2] C. ARMEN AND C. STEIN, *A 223 superstring approximation algorithm*, Discrete Applied Mathematics, 88 (1998), pp. 29 – 57. Computational Molecular Biology {DAM} - {CMB} Series.
- [3] G. AUSIELLO, P. CRESCENZI, V. KANN, MARCHETTI-SP, G. GAMBOSI, AND A. M. SPAC-CAMELA, *Complexity and Approximation: Combinatorial Optimization Problems and Their Approximability Properties*, Springer, Jan. 2000.
- [4] A. BLUM, T. JIANG, M. LI, J. TROMP, AND M. YANNAKAKIS, *Linear approximation of shortest superstrings*, J. ACM, 41 (1994), pp. 630–647.
- [5] J. GALLANT, D. MAIER, AND J. STORER, *On finding minimal length superstrings*, Journal of Computer and System Sciences, 20 (1980), pp. 50 – 58.
- [6] M. R. GAREY AND D. S. JOHNSON, *Computers and Intractability: A Guide to the Theory of NP-Completeness*, W. H. Freeman & Co., New York, NY, USA, 1979.
- [7] B. MA, *Why greed works for shortest common superstring problem*, Theoretical Computer Science, 410 (2009), pp. 5374 – 5381. Combinatorial Pattern Matching.
- [8] D. MAIER AND J. STORER, *A note on the complexity of the superstring problem*, Tech. Rep. 223, Department of Electrical Engineering and Computer Science, Princeton University, Princeton, New Jersey, 1977.
- [9] M. MIDDENDORF, *More on the complexity of common superstring and supersequence problems*, Theoretical Computer Science, 125 (1994), pp. 205 – 228.
- [10] K. PLOCIENNIK, *A probabilistic {PTAS} for shortest common superstring*, Theoretical Computer Science, 522 (2014), pp. 44 – 53.
- [11] J. S. TURNER, *Approximation algorithms for the shortest common superstring problem*, Information and Computation, 83 (1989), pp. 1 – 20.
- [12] A. ZARITSKY AND M. SIPPER, *Coevolving solutions to the shortest common superstring problem*, Biosystems, 76 (2004), pp. 209 – 216. Papers presented at the Fifth International Workshop on Information Processing in Cells and Tissues.

VITA

JR is BS Information Technology student of the Department of Computer Science at the Ateneo de Naga University.