

# Investigating Stochastic Diffusion Search in DNA sequence assembly problem

Fatimah Majid al-Rifaie  
Department of Computing  
Goldsmiths, University of London  
International Programme  
New Cross, London SE14 6NW  
Email: fmar2@student.london.ac.uk

Mohammad Majid al-Rifaie  
Department of Computing  
Goldsmiths, University of London  
New Cross, London SE14 6NW  
Fax: +44 (0)20 7919 7853  
Email: m.majid@gold.ac.uk

**Abstract**—This paper introduces a novel study on the performance of Stochastic Diffusion Search (SDS) – swarm intelligence algorithm – to address DNA sequence assembly problem. This is an NP-hard problem and one of the primary problems in computational molecular biology that requires optimisation methodologies to reconstruct the original DNA sequence. In this work, SDS algorithm is adapted for this purpose and several experiments are run in order to evaluate the performance of the presented technique over several frequently used benchmarks. Given the promising results of the newly proposed algorithm and its success in assembling the input fragments, its behaviour is further analysed, thus shedding light on the process through which the algorithm conducts the task. Additionally the performance of the algorithm is compared against several other techniques, demonstrating its weaknesses and strength in the experiments presented in the paper.

**Keywords**—DNA sequence assembly; Stochastic Diffusion Search; Multi-agent algorithm; Swarm intelligence

## I. INTRODUCTION

Every single cell in the body has a complete copy of about 3.2 billion<sup>1</sup> DNA base pairs or letters which build the human genome [1]. DNA has all the information necessary to build the whole living organism. Although the letters of the genetic alphabet Adenine (A), Thymine (T), Cytosine (C) and Guanine (G) are meaningless on their own, they are joined into useful instructions in genes. It is interesting to note that more than 99 percent of human's structure is genetically identical [2].

Imagine having several copies of the same book written in a language you cannot understand. Every page of each copy has been randomly cut into horizontal strip and a piece from one copy may overlap a piece from another copy. Assuming that some of strips are missing and some are splashed with ink, and maybe some of the books have random typos and error throughout, in different places. Try to arrange all the strips and assemble a single copy of the original book without any typos or errors. This process is similar to the important task of DNA sequencing.

In this work, a novel application of a swarm intelligence technique is introduced as a proof of principle. The swarm intelligence algorithm used is Stochastic Diffusion Search (SDS) which has a good potential to work in large search

spaces and noisy environments. This algorithm is explained in the paper and its application to the problem is detailed.

This paper starts by presenting the swarm intelligence algorithm along with a simple example demonstrating its use. Then a brief introduction is given to the DNA assembly problem and the solutions offered so far using swarm intelligence techniques. Subsequently, some experiments are designed and the performance of SDS is investigated using various benchmarks and then its performance is contrasted against several other techniques. Finally, the reason behind using SDS is further elaborated and the difference between utilising SDS and Smith-Waterman algorithm is discussed, followed by a conclusion and directions for future research.

## II. SWARM INTELLIGENCE

The paper is based on swarm intelligence which is one of the categories of artificial intelligence. Swarm intelligence is based on the study of behaviour of simple individuals (e.g. ant colonies, bird flocking, and honey bees, animal herding) that mimics the behaviour of swarms of social insects or animals [3]. More and more researches are interested in this field as swarm intelligence offers new ways of designing intelligence systems.

Among the successful examples of optimisation techniques inspired by swarm intelligence are: ant colony optimisation (inspired by foraging behaviour of real ant colonies) and particle swarm optimisation (inspired by bird flocking) [3]. In this work, Stochastic Diffusion Search (SDS) [4], [5] algorithm is used. This algorithm also belongs to the category of swarm intelligence and is based on mimicking the foraging behaviour of one type of ants *Leptothorax acervorum*. More details about SDS are provided in Algorithm 1 and a simple example is presented next.

### A. Search example with SDS

In the following example the aim is to find a 4-letter model (Table I) in a 32-letter search space (Table II).

There are four agents; and a hypothesis identifies four adjacent letters in the search space (e.g. hypothesis '6' refers to D-N-A-F; hypothesis '17' refers to A-S-S-E, etc.). In the first step, each agent initially picks a random hypothesis from the search space (see Table III). Assume that:

<sup>1</sup>In American English, 1 billion is equated to a thousand million (i.e. 1,000,000,000).

TABLE I: MODEL

Index:	0	1	2	3
Model:	D	N	A	F

TABLE II: Search Space

Index:	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
Search space	T	H	I	S	I	S	D	N	A	F	R	A	G	M	E	N

Index:	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31
Search space	T	A	S	S	E	M	B	L	Y	P	R	O	B	L	E	M

**Algorithm 1** SDS algorithm

**Initialisation phase:** Allocate agents to random hypotheses in the search space  
Until (all agents congregate on the best hypothesis)

- **Test phase**
  - Each agent evaluates its hypothesis
  - Each agent is classified into active or inactive
- **Diffusion phase**
  - Each inactive agent randomly chooses another agent to communicate with. If the inactive agent selects another inactive agent, no information will be transferred between the agents. Therefore the selecting agent should choose another hypothesis randomly. If the selected agent is active, the active agent communicate its hypothesis to the selecting agent

End

- The first agent points to the 27<sup>th</sup> entry of the search space and randomly picks one of the letters (e.g. the fourth one, (B)): 

O	B	L	E
---	---	---	---
- The second agent points to the 14<sup>th</sup> entry and randomly picks the first letter (E): 

E	N	T	A
---	---	---	---
- The third agent refers to the 8<sup>th</sup> entry in the search space and randomly picks the second letter (F): 

A	F	R	A
---	---	---	---
- The fourth agent goes the 20<sup>th</sup> entry and randomly picks the third letter (B): 

E	M	B	L
---	---	---	---
- The fifth agent refers to the 4<sup>th</sup> entry in the search space and randomly picks the second letter (S): 

I	S	D	N
---	---	---	---

The letters picked are compared to the corresponding letters in the model that is D-N-A-F (see Table I). In this case:

- The fourth letter from the first agent (E) is compared against the fourth letter from the model (F) and because they are not the same, the agent is set inactive.
- For the second agent, the first letter (E) is compared with the first letter from the model (D) and because they are not the same, the agent is set inactive.

TABLE III: INITIALISATION AND ITERATION 1

Agent No:	1	2	3	4	5
Hypothesis position	27	14	8	20	4
	OBLE	ENTA	AFRA	EMBL	ISDN
Letter picked:	4 <sup>th</sup>	1 <sup>st</sup>	2 <sup>nd</sup>	3 <sup>rd</sup>	2 <sup>nd</sup>
Status:	×	×	×	×	×

- For the third, fourth and fifth agents, letters ‘F’, ‘B’ and ‘S’ are compared against ‘N’, ‘A’ and ‘N’ from the model. Since none of the letters correspond to the letters in the model, the status of the agents are set inactive.

In the next step, each inactive agent chooses another agent and gets the same hypothesis if the selected agent is active. If the selected agent is inactive, the choosing agent generates a random hypothesis. Assume that the first agent selects the third one; since the third agent is inactive, the first agent chooses a new random hypothesis from the search space (e.g. 6). Fig.1 shows communication between agents.

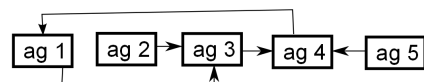


Fig. 1: Agent Communication 1.

The process is repeated for the other four agents. When the agents are inactive, they all choose new random hypotheses (see Table IV).

TABLE IV: ITERATION 2

Agent No:	1	2	3	4	5
Hypothesis position	1	6	22	12	17
	HISI	DNAF	BLYP	GMEN	ASSE
Letter picked:	1 <sup>st</sup>	2 <sup>nd</sup>	4 <sup>th</sup>	3 <sup>rd</sup>	3 <sup>rd</sup>
Status:	×	√	×	×	×

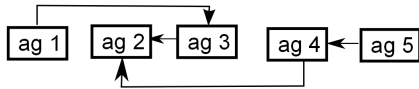


Fig. 2: Agents Communication 2.

In Table IV, the first, third, fourth and fifth agents do not refer to their corresponding letter in the model, therefore they become inactive. The second agent, with hypothesis ‘6’, chooses the second letter (N) and compares it with the second letter of the model (N). Since the letters are the same, the agent becomes active.

In this case, consider the following communication between the agents: (see Fig. 2)

- The third and fourth agents choose the second one
- The first agent chooses the third one
- The fifth agent chooses the fourth one

At this stage, the first and fifth agents, which chose the inactive third and fourth agents, have to choose other random hypotheses from the search space. However, agents three and four use the hypothesis of the active agent, two.

This process is repeated until all agents are active pointing to the location of the model inside the search. Depending on the problem, there are alternative termination strategies; for instance, in some cases, SDS algorithm is set to terminate only if all agents are active and refer to the same hypothesis.

The next section, provides a brief introduction to DNA assembly problem, stating the main phases and the major challenges faced by researchers in this field. This is followed by an overview of some of the algorithms that aimed to address the problem. Afterwards the experiments and results are reported.

### III. UNDERSTANDING DNA ASSEMBLY

There is no single solution available for NP-hard problems [1] and it is often not possible to find an extremely good algorithm that solves such problems [6].

In DNA assembly, a process is required to join the relevant fragments together. In other words, the overlapping fragments are to be assembled back into the original DNA sequence. Therefore, the goal of genome projects is to reconstruct the original genome sequence of an organism. To achieve the goal, DNA fragment assembly process is divided into three phases [7], [8]:

- 1) **Overlap Phase** is tasked to find the common sequence among the prefix of one sequence and suffix of another.
- 2) **Layout Phase** uses alignment strategies to determine the order of fragments based on high overlap scores and according to the level of similarity.
- 3) **Consensus Phase** assembles all fragments into the consensus sequence and omits the similar parts.

The quality of a consensus sequence is measured by the term coverage [9], [7]. Coverage is evaluated according to the following equation:

$$\text{Coverage} = \frac{\sum_{i=1}^n \text{length of fragment } i}{\text{target sequence length}} \quad (1)$$

where  $n$  is the number of fragments.

The higher the coverage, the higher the probability of covering original genome, the higher the correctness of the assembled parts, the fewer the number of the gaps, and the better the result [7], [10].

The Layout Phase is the most complex step due to the difficulty of finding the best overlap. This difficulty is caused by the following challenges [10], [7]:

- **Unknown orientation:** After the original sequence is divided into many fragments, the direction may change.
- **Base call errors:** substitution, insertion, and deletion errors are types of base call error. The errors happen because of experimental errors in the electrophoresis procedure that affects the finding of fragment overlaps.
- **Incomplete coverage:** It occurs when the algorithm cannot assemble a given fragments into one contig.
- **Repeated regions:** the problem occurs when some sequences are repeated two or more times in the DNA. None of the current assembly programs can solve the problem without an error [6].
- **Chimeras and contamination:** Chimeras arise when two fragments that are not adjacent, or overlapping on the target molecule, join together into one fragment. Contamination occurs due to the incomplete purification of the fragment from the vector DNA.

#### A. DNA Sequence Assembly and Swarm Intelligence

DNA Assembly problem is still open to a large extent because of the principal issue of “scaling up to real organism”. Some of the swarm intelligence and evolutionary algorithms, such as genetic algorithms and ant colony optimisation have been used for the fragment assembly problem focusing on the overlap, layout and consensus approach [11].

In 1995, Rebecca Parsons and Johnson created performance improvements for a genetic algorithm applied to the DNA sequence assembly problem [12]. In 2003 Kim and Mohan used a new parallel hierarchical adaptive genetic algorithm. The method is reported as accurate and noise-tolerant compared to previous methods [13]. In the same year, Meksangsouy and Chaiyaratana proposed ant colony optimisation. The goal of the search was to find the right order and orientation of each fragment to create a consensus sequence [14]. In 2005 Fang proposed approach speeded up the searching process and maximised the similarity or overlaps between given fragments [15]. Alba and Luque presented several methods, including genetic algorithm, a CHC method, scatter search algorithm, and simulated annealing to solve accurately DNA Assembly problem in 2005 [16]. They also proposed a local search

method named PALS in 2007 [17]. In 2008, Luque and Alba studied the behaviour of a hybrid heuristic algorithm that combines a heuristic, PALS, with a meta-heuristic, a genetic algorithm, achieving an assembler to find optimal solutions for large instances of this DNA assembly problem [18]. In 2010 Kubalik presented a method called Prototype Optimisation with Evolved Improvement Steps (POEMS). Also in the same year Minetti and Alba presented a paper about how noiseless and noisy instances of this problem are handled by three algorithms: problem aware local search, simulated annealing and genetic algorithms [19].

There are some other solutions that are proposed in 2011 for DNA sequence assembly problem using Particle Swarm Optimisation (PSO) with Shortest Position Value (SPV) rule [1]. In 2012 Firoz analysed and discussed the performance of two swarm intelligence based algorithms namely Artificial Bee Colony (ABC), and Queen Bee Evolution Based on Genetic Algorithm (QEGA) to solve the fragment assembly problem [20]. In 2013 Fernandez-Anaya et al. designed a nature inspired algorithm (PPSO+DE) based on Particle Swarm Optimisation and Differential Evolution [21].

The next section, presents the experiments designed for this paper, demonstrating the performance of the proposed algorithm. Then the results are reported along with comparisons against other techniques.

#### IV. EXPERIMENTS AND RESULTS

In order to understand the process through which SDS is adopted and adapted for DNA sequence assembly problem, a number of fragments are used in the experiments. The fragments are the input of the program and the program is responsible to assemble the fragments and create one long sequence. This is achieved by taking a fixed number of characters from the end of the first fragments and trying to find those characters in the other fragments using SDS algorithm. Once the other fragment is found, the two fragments are joined and the repeated part is deleted from one of the fragments. This will create a longer fragment. This process is repeated until all fragments are joined. The steps required for SDS to assemble a set of fragments are detailed in Algorithm 2.

---

#### Algorithm 2 DNA sequence assembly using SDS

---

Choose a model from the end of 1<sup>st</sup> fragment in the search space

While (true)

- Use SDS to search the model in the fragments
  - If no matching fragment is found
    - Choose the model from the beginning of the first fragment
    - Use SDS to search the model in the fragments
    - If no matching fragment is found  
Break
- Compile a list of fragments where the model is found
- Pick the fragment ( $j^{th}$ ) with the maximum similarity (based on agents activity)
  - assemble fragments  $i$  and  $j$ .
  - Delete the  $j^{th}$  fragment
  - Choose a new model from the end of assembled fragment

End While

---

In the experiments reported in this paper the agent size is empirically set to 100 and the model size for SDS is set to 50.

TABLE V: Benchmark datasets

Benchmark	Coverage	Mean fragment length	Number of fragments	Original sequence length
x60189 4	4	395	39	
x60189 5	5	286	48	3,835
x60189 6	6	343	66	
x60189 7	7	387	68	
m15421 5	5	398	127	
m15421 6	6	350	173	10,089
m15421 7	7	383	177	
j02459 7	7	405	352	20,000
bx842596 7	7	703	773	77,292

Table V, as proposed by Mallén-Fullerton et al. [11], shows the benchmarks used by SDS for DNA assembly.

Using the benchmarks provided, SDS algorithm assembles the entire sequences correctly. Table VI shows the performance of SDS when assembling the nine aforementioned benchmarks. Each benchmark is assembled 50 times. As the table shows, the larger the coverage, the more SDS iterations it takes to fully assemble the datasets. While the number of overall algorithm cycles needed follow the same structure, there are some exception caused by the order of the fragments. Observing the sum of active agents over all the iterations and their consistent proximity (check the negligible difference between the median and the mean, as well as the value of the standard deviation) shows the robustness of the technique.

The results shown in Table VI indicate that three of the benchmarks ( $m15421$  6,  $m15421$  7 and  $bx842596$  7) are not assembled fully into one sequence. SDS has been able to assemble two large, accurate sequences from the fragments of each of these datasets which make up the whole dataset. However up to this point, given there were no similarities between the two resulting sequences, they are returned separately. Therefore, caution is taken and they are reported as not completely assembled.

Next, one of the benchmarks is chosen ( $x60189$  4) and the analysis are reported based on this benchmark. The results are compatible with the ones generated from the other benchmarks. Fig. 3-left shows the level of agents activity at various stages of SDS assembling process, including both when a match is found and when a match is not found in any given fragment. The activity of the agents is in the range  $[0, 100]$ , however if less than the entire agent population (i.e. 100) are active, the agents' hypotheses are not taken into account for the assembling purpose; this ensures the presence of a full match. Reducing the 100% accuracy would cater for a noisy environment which is one of the strengths of SDS algorithm.

To provide a better understanding, the histogram of the agents' activity is presented in Fig. 3-middle. This graph clearly shows that in most cases there is no high similarity between fragments (note that the similarity between fragments is evaluated by comparing the model to the fragments). How-

TABLE VI: Summary of Assembling Four Datasets

	Cycles	SDS Itrs	Sum of active agents				
			Median	Mean	Stdev	Max	Min
<b>x60189 4</b>	23	46,899	384,075	384,746	4,007	392,361	374,717
<b>x60189 5</b>	17	53,649	418,777	418,744	4,119	430,935	409,683
<b>x60189 6</b>	28	114,799	752,539	752,176	4,382	763,060	740,365
<b>x60189 7</b>	26	124,249	1,095,673	1,096,029	5,928	1,109,507	1,083,767
<b>m15421 5</b>	57	476,149	2,175,028	2,176,785	9,722	2,220,938	2,150,915
<b>m15421 6</b>	-	-	-	-	-	-	-
<b>m15421 7</b>	-	-	-	-	-	-	-
<b>j02459 7</b>	129	3,174,449	12,092,968	12,087,911	21,613	12,127,170	12,021,776
<b>bx842596 7</b>	-	-	-	-	-	-	-

ever when there is a match (i.e. 100% activity), the fragments are joined on the fly.

Fig. 3-right provides a close-up view of the graph on its left and demonstrates that when there are no exact matches, some of the SDS agents could be activated; however if there are no full match, the activated agents eventually lose their active status in the consequent iterations when they choose a different micro-feature. This feature is particularly useful in a noisy environment whose complete analysis will be provided in an expanded future publication.

In a similar experiment and in order to analyse the behaviour of the agents when some of the fragments are contaminated with noise, some noise (i.e. of type substitution) is added to all the fragments. Fig. 4-left shows the activity of the agents and Fig. 4-middle and right illustrate the frequency of activity level at various iterations. Note that there are fewer number of iterations needed before SDS terminates (as at some point during the process, no match is found from either end of the growing sequence). However the proportion of agents activity between 0 and 100 is increasing with the presence of noise. Despite the fact that the entire fragment is contaminated with noise, SDS is able to produce more than half the length of the target sequence. Further research is required to improve this rate.

In order to illustrate the activity of the agents at each SDS iterations, the graphs in Fig. 5 are presented. In Fig. 5-left the activity of SDS agents are displayed when there is no match. As shown, most of the agents are inactive and very a few flicker from being active and then back to being inactive.

However, on the contrary to the lack of a match, when there is a full match, as shown in Fig. 5-right, soon after the start of the SDS iteration and through agents' communication and information exchange, the entire population becomes active and points to the right position, which is the position of the model within the fragment.

#### A. Comparison with other techniques

In another analysis, the performance of SDS is compared against a few other algorithms tasked with assembling the benchmarks. These algorithms, which are used in this context

TABLE VII: Comparison with other techniques

	SDS	PALS	GA	PMA	CAPS	Phrap
<b>x60189 4</b>	1	1	1	1	1	1
<b>x60189 5</b>	1	1	1	1	1	1
<b>x60189 6</b>	1	1	-	1	1	1
<b>x60189 7</b>	1	1	1	1	1	1
<b>m15421 5</b>	1	1	6	1	2	1
<b>m15421 6</b>	2	NA	NA	NA	NA	NA
<b>m15421 7</b>	2	1	1	2	2	2
<b>j02459 7</b>	1	1	13	1	1	1
<b>bx842596 7</b>	2	2	-	2	2	2

in the literature, are genetic algorithm (GA), a pattern matching algorithm (PMA), Problem Aware Local Search (PALS) and commercially available packages: CAP3 and Phrap. The algorithms are compared in terms of the final number of contigs assembled. Despite being in the early stages of its application in DNA assembly problem, SDS shows a competitive performance (see Table VII<sup>2</sup>). Other than an isolated case (m15421 7), where PALS and GA outperform SDS, in the rest of the cases (89%), SDS either presents similar or better outcome. SDS is also tried on a benchmark (m15421 6) that is not attempted by the rest of the techniques. The accuracy of the assembled sequences is 100%; in other words, whenever the accuracy is less than 100%, the results are considered unsuccessful. In these experiments, SDS deals with various issues common in DNA sequence assembly<sup>3</sup>, including but not limited to fragments with varying lengths, unknown orientation, incomplete coverage, repeated regions, chimeras and contamination, etc.

#### B. SDS vs. Smith-Waterman algorithm

Many DNA sequence assembly techniques use Smith-Waterman algorithm [23], which is a pairwise alignment method to create a similarity matrix between the fragments, therefore generating a complete picture of the entire available

<sup>2</sup>The results of these algorithms, other than SDS, are borrowed from [22].

<sup>3</sup>These issues are explained in section III.

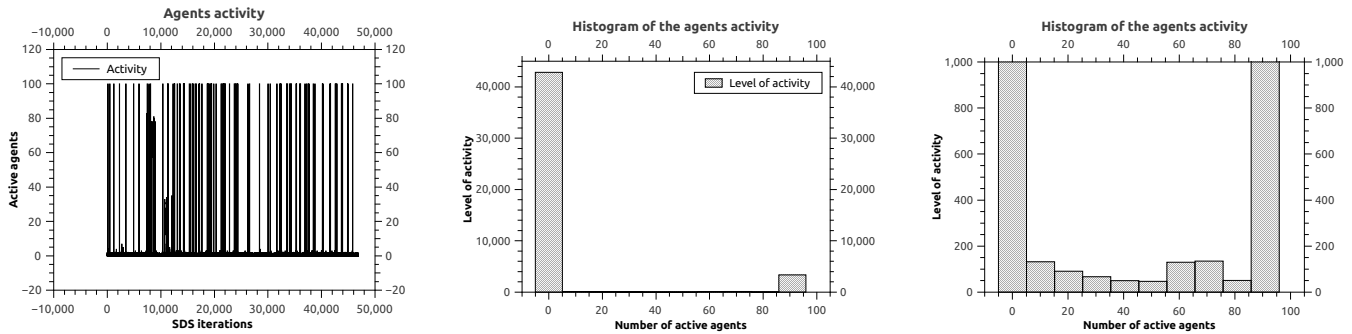


Fig. 3: Left: activity of the agents in the fragments of  $x60189\ 4$ ; middle: the histogram of the activity of the agents; right: zooming to show the activity of agents between 0 and 100.



Fig. 4: Left: activity of the agents in noisy set of fragments; middle: histogram of the activity of the agents in noisy set of fragments; right: zooming to show the activity of agents between 0 and 100.

data before setting off to the overlapping and assembling stage. While Smith-Waterman algorithm provides a precise and detailed account of the input data, it comes at the expense of being time consuming and computationally expensive [24].

Assuming there are  $n$  fragments, once the similarity between each pair is calculated using Smith-Waterman algorithm, an  $n \times n$  matrix is created. The matrix is then used by other optimising algorithm to conduct the overlapping phase. The results of many of these algorithms are reported in [11].

In the experiments reported earlier in the paper, instead of using Smith-Waterman algorithm to calculate the similarities between fragments, SDS picks a model from a given fragment and aims to find the model in the rest of the fragments. Among the fragments containing the model, the one with the highest similarity is picked and assembled on the fly and then removed from the search space, thus reducing the subsequent computational cost.

On the contrary to many other swarm intelligence and evolutionary computation, SDS has been successful in assembling the benchmarks without using Smith-Waterman algorithm, therefore avoiding its time consuming and computational expensive nature. To understand the full picture of the process, further analysis is needed, among other things, to verify the impact left on the assembling process without accessing the very detailed information provided by Smith-Waterman algorithm.

## V. CONCLUSIONS

Since DNA fragment assembly problem is NP-hard, it is difficult to find optimal solutions. The increasing presence of biological data and the requirements to study and understand them closely lead us to use computational approaches. In addition to a brief literature review given in this work, it is shown that DNA fragment assembly problem can be attempted with meta-heuristics.

For the first time, this paper presents Stochastic Diffusion Search (SDS), which belongs to the extended family of swarm intelligence algorithms, in the context of DNA fragment assembly problem. An initial study into the behaviour of SDS is provided, offering an analysis into the agents' activity using several benchmarks. The results are promising as they demonstrate how the activity of the agents shed light into the way agents interact and eventually finalise the assembling process. Additionally it is shown that the level of agents' activity provides a measure of similarity between fragments, thus allowing more similar fragments to be joined in the assembling process.

As part of the future research, this algorithm will be compared against other evolutionary computation techniques used in this field; also larger datasets with more complex features are to be used, and more research is needed in order to theoretically determine the two values (population size and model size) of the SDS parameters. Additionally, CPU time and RAM memory usage will be taken into account for all the

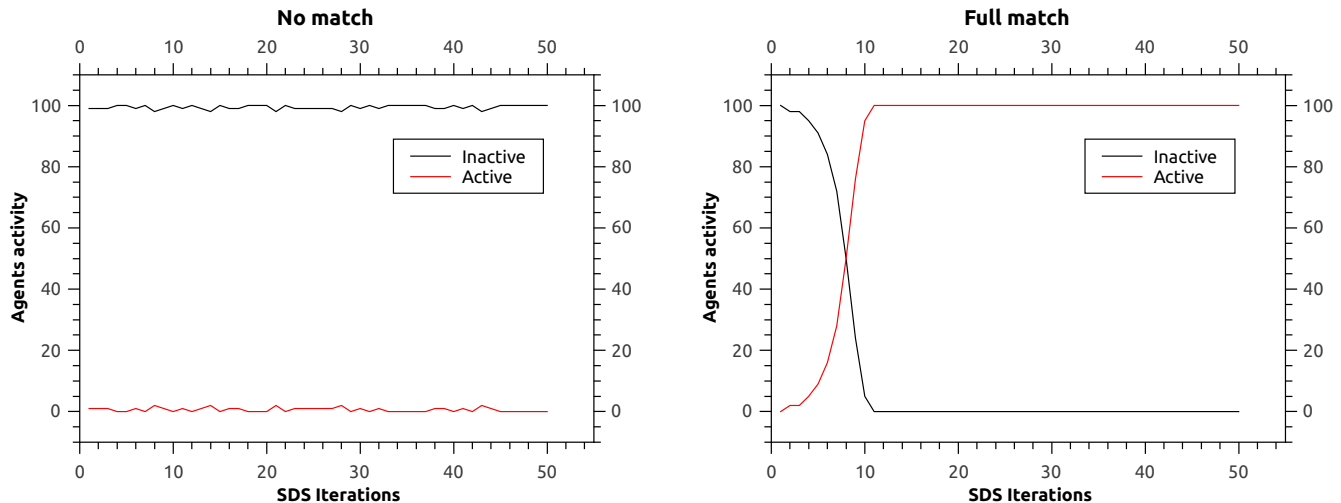


Fig. 5: Left: absence of a match; right: presence of a full match.

comparisons to provide a more comprehensive account on the performance of the proposed algorithm.

#### REFERENCES

- [1] R. S. Verma, V. Singh, and S. Kumar, "Dna sequence assembly using particle swarm optimization." *International Journal of Computer Applications*, vol. 28, 2011.
- [2] U.S. National Library of Medicine, "Cells and dna. what is dna?" <http://ghr.nlm.nih.gov/handbook/basics/dna>, accessed: 06/01/2015.
- [3] C. Blum and X. Li, *Swarm intelligence in optimization*. Springer, 2008.
- [4] M. M. al-Rifaie and M. Bishop, "Stochastic diffusion search review," in *Paladyn, Journal of Behavioral Robotics*. Paladyn, Journal of Behavioral Robotics, 2013, vol. 4(3), pp. 155–173.
- [5] J. Bishop, "Stochastic searching networks," in *Proc. 1st IEE Conf. on Artificial Neural Networks*, London, UK, 1989, pp. 329–331.
- [6] P. Pevzner, *Computational molecular biology: an algorithmic approach*. MIT press, 2000.
- [7] C. Cotta, A. Fernández, J. Gallardo, G. Luque, and E. Alba, "Metaheuristics in bioinformatics: Dna sequencing and reconstruction," *Optimization Techniques for Solving Complex Problems*, pp. 265–286, 2009.
- [8] K.-W. Huang, J.-L. Chen, C.-S. Yang, and C.-W. Tsai, "A memetic particle swarm optimization algorithm for solving the dna fragment assembly problem," *Neural Computing and Applications*, pp. 1–12, 2014.
- [9] J. C. Setubal, J. Meidanis, and J. C. Setubal-Meidanis, *Introduction to computational molecular biology*. PWS Pub., 1997.
- [10] L. Li and S. Khuri, "A comparison of dna fragment assembly algorithms." in *METMBS*, vol. 4, 2004, pp. 329–335.
- [11] G. M. Mallén-Fullerton, J. A. Hughes, S. Houghten, and G. Fernández-Anaya, "Benchmark datasets for the dna fragment assembly problem," *International Journal of Bio-Inspired Computation*, vol. 5, no. 6, pp. 384–394, 2013.
- [12] R. Parsons and M. E. Johnson, "Dna sequence assembly and genetic algorithms-new results and puzzling insights." in *ISMB*, 1995, pp. 277–284.
- [13] K. Kim and C. K. Mohan, "Parallel hierarchical adaptive genetic algorithm for fragment assembly," in *Evolutionary Computation, 2003. CEC'03. The 2003 Congress on*, vol. 1. IEEE, 2003, pp. 600–607.
- [14] P. Meksangsouy and N. Chaiyaratana, "Dna fragment assembly using an ant colony system algorithm," in *Evolutionary Computation, 2003. CEC'03. The 2003 Congress on*, vol. 3. IEEE, 2003, pp. 1756–1763.
- [15] S.-C. Fang, Y. Wang, and J. Zhong, "A genetic algorithm approach to solving dna fragment assembly problem," *Journal of Computational and Theoretical Nanoscience*, vol. 2, no. 4, pp. 499–505, 2005.
- [16] G. Luque and E. Alba, "Metaheuristics for the dna fragment assembly problem," *International Journal of Computational Intelligence Research*, vol. 1, 2005.
- [17] E. Alba and G. Luque, "A new local search algorithm for the dna fragment assembly problem," in *Evolutionary Computation in Combinatorial Optimization*. Springer, 2007, pp. 1–12.
- [18] —, "A hybrid genetic algorithm for the dna fragment assembly problem," in *Recent Advances in Evolutionary Computation for Combinatorial Optimization*. Springer, 2008, pp. 101–112.
- [19] G. Minetti and E. Alba, "Metaheuristic assemblers of dna strands: Noiseless and noisy cases," in *Evolutionary Computation (CEC), 2010 IEEE Congress on*. IEEE, 2010, pp. 1–8.
- [20] J. S. Firoz, M. S. Rahman, and T. K. Saha, "Bee algorithms for solving dna fragment assembly problem with noisy and noiseless data," in *Proceedings of the fourteenth international conference on Genetic and evolutionary computation conference*. ACM, 2012, pp. 201–208.
- [21] G. M. Mallén-Fullerton and G. Fernández-Anaya, "Dna fragment assembly using optimization," in *Evolutionary Computation (CEC), 2013 IEEE Congress on*. IEEE, 2013, pp. 1570–1577.
- [22] E. Alba and G. Luque, "A new local search algorithm for the dna fragment assembly problem," in *Evolutionary Computation in Combinatorial Optimization*. Springer, 2007, pp. 1–12.
- [23] T. F. Smith and M. S. Waterman, "Identification of common molecular subsequences," *Journal of molecular biology*, vol. 147, no. 1, pp. 195–197, 1981.
- [24] G. Wieds, "Bioinformatics explained: Blast versus smith-waterman," *CLCBio*. <http://www.clcbio.com/index.php>, 2007.