

FORENSIC BIOMETRICS IDENTIFICATION SYSTEM FOR DNA PROFILE HUMAN BASED ON ASSOCIATION RULES

Najah H. Faleh*¹, Karim H. Al-Saedi²

^{1&2} *Department of Computer Science, College of Science, Mustansiriyah University, Baghdad, Iraq.*

najah.19@yahoo.com

Abstract It should be noted that there is a joint work between biomedicine such as forensic medicine and information technology through the use of information technology technologies in all fields, including determining the of DNA profile. One of the leading biotechnologies in this field is data mining techniques. There are many ways to identify disaster victims, such as fingerprints, dental record and DNA profile matching. DNA matching is a highly accurate identification way that does not need specific parts of the victim's body. Deoxyribose Nucleic Acid (DNA) is the basic elements that make up an entire section of a human. The core elements store unique information for each individual and will be passed on through generations. DNA also helps in identifying the father in paternity testing,. The limitation of applying DNA matching for disaster victim identification lies on expensive and time consuming process. To address this situation, in this paper, we performed a method to measure the confidence of matching of human DNA profiles identification using Association Rule Classification System is proposed. In this Classification system, DNA profile data is used as an input that stores human identity along with its DNA profile. advisable information or good patterns from present datasets for certain objective. The results were satisfactory and characterized by large percentage and high accuracy. Finally, performance of this system is evaluated and in turns the proposed system proves its capability in forensic human identification and scalability to handle huge amounts of data.



Keywords: *Association Rule Classification, Biometric system, Data mining, DNA-profile, SQL server.*

1. INTRODUCTION

Data mining as Association Rules Classification have established a plenty of effective applications in DNA profile analysis and human identification [1]. Data mining is a fast-growing subarea of machine learning that handle discovering meaningful knowledge from large datasets [2]. Nowadays, DNA motifs technology has formed a new line of research in both data mining and bioinformatics [1]. DNA is used to identify criminals, to verify suspects, and to innocent persons accused or wrongly convicted of crimes, with incredible accuracy when there is biological evidence [3]. The DNA profile is related to the prevention or detection of crimes related to the identification of persons in the National Security Service or in an anti-terrorism investigation [4]. If there is damage to all or part of the victim's or suspect's body, there is little evidence of crime, there will be difficulties in identifying and this causes problems, and issues such as the length of the case's settlement[5]. Therefore, victims or suspects' DNA is used as the primary means of identification. This is done because DNA can be found in almost every human body, as well as the unique properties of DNA that can be used for identification. DNA identification consists of sample analysis to isolate a unique set of DNA markers [6]. The analyst then compares the DNA profiles to determine

whether the person's DNA sample matches the evidence obtained from the crime scene or from a family relationship [7]. In a family, a child's DNA profile is a combination of the DNA profiles of both parents because the child has one allele that he inherits from his father and the other from the mother. An individual can be declared a child of a father or mother if he has a similar DNA that is about 50% of the total DNA profile, because 50% of the DNA is inherited directly by the father or mother [8]. Moreover, it can be concluded that matching with a biological novelty or ancestor is about 25%, as well as between the individual and his siblings (full sibling) from 45% to 54%. While the ratio between the individual and his half-brother (half-siblings) is 25%. Likewise, for an uncle or aunt, he is also similar to 25% and 25% with a nephew who has the same parental relationship. From there it can be concluded that there is only 12.5% similarity to cousins [9]. DNA profiles will be compared have 16 locus, any DNA loci is involve of two alleles, of the marker-one inherited from the mother and one inherited from the father. Sixteen locus should be compared everything to make a decision whether there is any relationship between the DNA profile evidence the biological by comparison [10]. In this paper we automate the usage of the physiological characteristics for identification and verification purpose by proposing a



biometrics system based on Association Rules Classification System as biometrics too [11]. Sixteen STRs loci for Population in forensic center in Baghdad are used and saved in a database. The database is built using a SQL Server environment. The famous biometric techniques shown in different research works include: fingerprint, iris, DNA, hand geometry, signature, voice... etc. [12]. The following Section 2 discusses the related works; Section 3 discusses the proposed Association Rule Classification System for DNA profile Identification. An experiment of the usage of Classification System in Section 4, followed by conclusions discussion in Section 5.

2. RELATED WORKS

Nurtami Soedarsono et al (2016), [13] proposed A Novel Human STR Similarity process utilize cascade statistical fuzzy rules with tribal information inference (NHSSM). This system suggest a new process for inferring the similarity of tribal knowledge utilize vague statistical rules to deal with uncertainty and inaccuracy in DNA profile, and to implicate tribal data inference on short tandem repeat-depend DNA similarity match utilize statistical fuzzy rules where the allele marker's statistical distribution probability density function is used as the membership function in the fuzzy rules of the initial FIS, the novel manner makes it potential to tell the tribal similarity among two STR profiles. The experience visible that the novel manner is capable to distinguish DNA typing among tribal groups.

Saja Dheyaa Khudhur (2017), [14] proposed the System Identify Human Based on DNA and Dental X-Ray as valuable and reliable biometrics tools (SIHBODAD). The motivation of this labor is the collection of the Dental features and DNA for individual identification. There are 16 locus used of (STRs) DNA profile and a bite-wing image, used to extractor the dental data, as DNA and dental X-Ray features profiles. The dental features are Standard Deviation (STD), Euler number and Intensity, extract from the bite-wing image through use a three phases algorithm. This phases are features extraction, image segmentation and classification. The purpose of the suggest system is to presented a forensic human identification system based on a create database.

Maria Susan Anggreainy et al (2019), [15] proposed, Family relation And STR-DNA Match utilizes fuzzy inference (FASM). The aim of this proposed system to diagnosis of the unknown person if the compare is between the victim and his

parents, or the unknown person parents have missing or far away from where the unknown person, it is essential to try to identify DNA profile to derivation of live family. It will be performed a certain manner in proposed system to measure the similarity of individual DNA typing utilize fuzzy similarity. Outcome in this proposed system is the value of person similarity with levels, specifically smaller, moderate and higher. It will be utilizing siblings as an exchange for a both parents, in this fuzzy system, in order to the value that creates it adequate and perfectly close to the value of comparison with both father and mother.

3. PROPOSED METHOD

To improve the accuracy of DNA analysis, the process of checking to do as much as possible so as to get objective results. This requires the presence of DNA samples in large numbers. The proposed Association Rule Classification System for DNA profile Identification as metrics to satisfy the identity of the human. It will be utilize a sixteen STR loci where the human genome contains a short tandem repeats refer to as STR. The adopted STR profiles are collected from the pre-stored STR loci of forensic center in Baghdad as 400 DNA profiles. The database consists of two tables: a source table and a reference table. In both tables there are 33 fields/columns namely Name (name of individual), id (serial number of DNA profile of individual), And 31 locus values of 16 locus of DNA profile (sixteen loci are CSF1PO, D13S317, D16S539, D18S51, D19S433, D21S11, D2S1338, D3S1358, D5S818, D7S720, D8S1179, FGA, TH01, TPOX, VWA and Amelogenin)[8]. The system offers a high ability in performing the database activities, such as insert, update and search the physiological characteristics (STR DNA profile for individuals). Moreover, the matching process is performed at a time of presence the query physiological characteristics with these the saved information. The software programs are utilized as Visual Studio C# and SQL server. The GUI provides the ordinary users a high flexibility and simplicity dealing with the system without need to pre-knowledge about the internal procedure. The proposed has the ability to determine an integrated human identification system based on Association Rules Classification System for STR DNA profile identification techniques, linked with a huge database. And also presents a description of designing and implementation phases of the proposed system, which comprises of four operations under two main phases which are: preprocessing phase and Classification phase as below. The figure (3.1) shows the Architecture of proposed system.

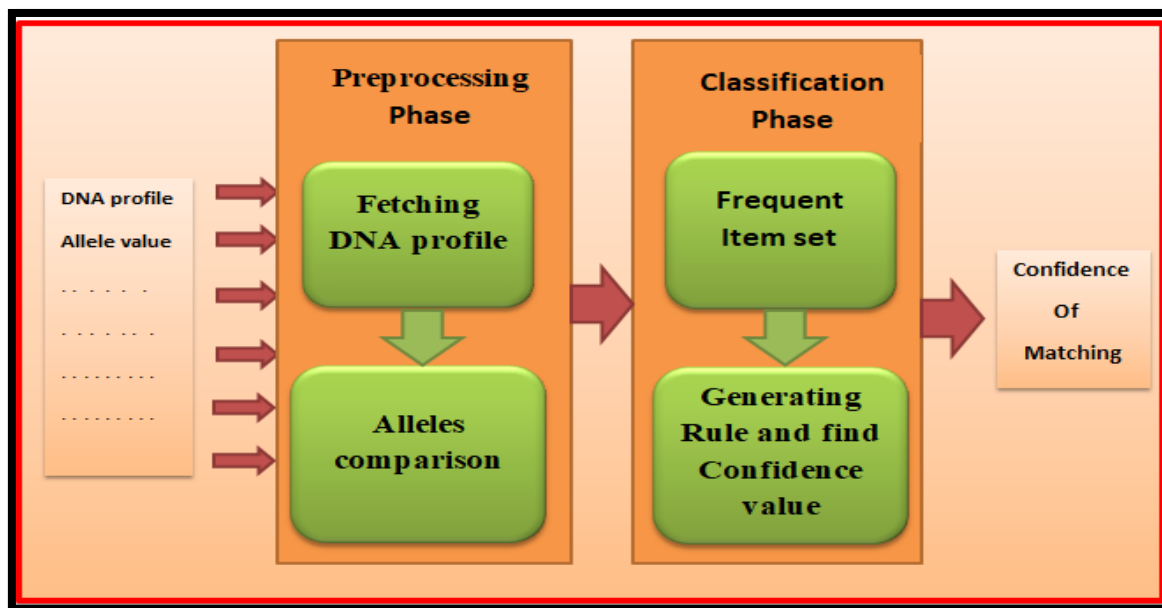


Figure (3-1) shows Association Rule Algorithm architecture.

3.1 Preprocessing.

This phase responsible for extracting the values of alleles in each loci from the STR DNA profile of the { source, reference} and to the management of comparison method. This phase has two main steps, namely Fetching DNA profile, and Allele comparison.

3.1.1 Fetching DNA Profile

Our data are stored in a database (SQL Server), in order to make the processing on data easier and scalable in data. In this step of this phase needs to fetch the data from database by opening a connection with SQL Server. If there is a connection then will execute a SQL statement which will fetch data from both source and reference table of STR DNA profile which will be ready in next step. Then Check the allele's fields are filled. If it is not filled, it must be refilled.

3.1.2 Alleles Comparison

The process of alleles filtration and arrangement between alleles. There are two alleles in each locus in DNA profile. Comparing the DNA profile of the victim (source) and the DNA profile of the reference person is done by comparing all alleles (Source allele_A with the Reference allele_A or Reference allele_B and Source allele_B with the Reference allele_A or Reference allele_B). The first allele group is from the DNA of a victim and the second allele group is from the DNA reference. Each DNA profile has 15 loci (30 alleles) and one amelogenin.

3.2 Classification Phase

This phase is detection of associations between items in huge transactional or relative data_sets. The detection of important Association Rule between large amounts of transaction records can assist in decision making process like determine the confidence of matching for DNA profile identification. This phase has two main steps as following:

3.2.1 Frequent Item Set Step

The itemset is a set of items associate together. when the itemset has k-items it will be namely a k-itemset. In this step uses Apriori algorithm being more convenience to the nature of the dataset, to frequent the k_itemsets from all itemset in the transaction to generate association rules. The frequent Item set is an item set which support value is maximal or equal than a threshold value (min_supp) which determine by user.

3.2.2 Generating Rule and find Confidence value step

This unite responsible for extracting the Association rules can be thought of as an IF-THEN rules relationship and compute confidence value. There are two elements of these rules:

- Antecedent (IF): This is an item (victim DNA profile) that is typically found in the Item sets or Datasets.
- Consequent (THEN): This (reference DNA profile) comes along as an item with an Antecedent.

This is where the Apriori Algorithm comes into play. Let's understand the math of these process. There are 2 ways to measure association:

- Support.
- Confidence.

Support: It gives away the fraction of transactions which consist victim DNA profile Alleles and reference DNA profile Alleles. Essentially support talk us about the frequently find items find frequently. So with this, it will be can filter out the items that have a low frequency. Where:

A: is a right hand side (victim DNA profile) .

B: is a left hand side (reference DNA profile).

N: is a count of all transactions, in our case, N equals 30, which is the number of alleles at 15 genetic locus, except amiloginine.

Confidence: It talk us how sometime the items (victim DNA profile) alleles and (reference DNA profile) alleles happen together, given the number times victim DNA profile occurs. Confidence = frequent (A , B) / frequent (A) (3.2) [16]
Where:

A: is a right hand side (victim DNA profile).

B: is a left hand side (reference DNA profile).

The algorithm of proposed system declare below in algorithm (3.1).

Algorithm 3.1 (Association Rules Classification System)

Input

DNA Profile (Alleles value)

Output

Confidence of matching

Begin

Step1: Prepare the connection string to SQL Server

*Step2: if there is a connection with database
then connect to database*

Else

message "There is error in connection"

Step 3: If there is a connection prepare SQL query to fetch DNA profile for Source and reference Alleles.

Step 4: Put the data in memory for next step in processing

*Step5: if ((Source(A) equal Reference(A) and Source(B) equal Reference(B)) or
((Source(A) not equal Reference(A) and Source(B) not equal Reference(B)) Then
Compar.(Source(A), Reference(A)) and (Source(B), Reference(B))*

Else

Compar.(Source(A), Reference(B)) and (Source(B), Reference(A))

Step6: Ck: Candidate item set of size k

Lk : frequent item set of size k

L1 = {frequent items};

for (k = 1; Lk != ∅; k ++) do begin

Ck+1 = candidates generated from Lk;

For each transaction t in database do

increment the count of all candidates in Ck+1 that are contained in t

Lk+1 = candidates in Ck+1 with min_support

If support >= min_supp. Then insert item into item set

Else

Remove item

End

Return U k Lk until no more frequent item sets found

4. EXPERIMENT

Many experiments were performed to assess the accuracy of the proposed system. Dataset utilized in the experiment is a STR DNA profile data acquired from the forensic medicine center in Baghdad. It consists of 200 DNA profile dataset as source sample and 200 DNA profile data as reference sample. The data is then saved in a database of DNA profiles for all individuals. For experimental similarity mensuration and measure the confidence of matching of victim DNA profile with the reference family relationship DNA profile. The data used contain of 400 sample DNA profile including data from

individuals with a biological relationship. As well as to gain a high accuracy of the results from the proposed system, will performed the process of identifying DNA profiles for individuals by applying the Association Rules Classification algorithm. To monitor the performance of DNA profile confidence of matching using ARC. Four comparison scenarios have been applied. The first scenario is to measure the confidence of matching between a child DNA profiles with a mother DNA profile. The second scenario is to measure the confidence of matching between a victim DNA profile and a son DNA profile. The third scenario is to measure the confidence of matching between two DNA

profiles that have no relation. The fourth scenario is the compare confidence of matching between the victims DNA profiles with itself within 400 DNA profiles. The comparison is made according to the proportions of family relationship, where each family relationship has a certain threshold, on the basis of that will be determine the confidence of matching and whether the victim has a family relationship with the reference or not. Table (4.1) show the percentage of Family relationship threshold [9] [15], which are will depending on it.

Table (4.1) Biological Family relationship Threshold

Biological Family relationship	Threshold
Same Person	(100 %)
Father , Mother , Son	(50 %)
grandmother , grandfather	(25 %)
Full_Sibling	(54 – 54 %)
Half_Sibling	(25 %)
Uncle, Aunt , Niece, Nephew	(25 %)
Cousin	(12.5 %)

4.1 First Scenario

The aim of such scenario is to know the accuracy and validate the confidence of matching of system by compare between two DNA profiles which have family relationship. Table (4.2) that shows both DNA profiles of individuals and confidence of matching value between child DNA profile and mother DNA profile using the proposed Association Rule Classification System (ARCS). The DNA profile matching result is 60 % it very acceptable. To verify the result, the comparison for this DNA set with threshold of family relationship. Where the child's DNA profile in family relationship is a mixture of the DNA profiles of both parents because the child has at each genetic loci the two alleles, one of them inherits from his father while the other inherits from his mother. Where it can tell and prove that the person belongs to a father or mother if he possesses at least 50% confidence of matching because it is the certain amount of the percentage that he inherits from his parents. The result produced by the proposed Association Rule Classification System is more realistic compared to if the result produced is lower than 50%. Likewise it was confirmed and verified

between the experimental and real ratios, where the conformity ratio was very high. Those two profiles have a strong family relationship, so the confidence of matching result should be high.

Table (4.2) First Scenario

Loci Name (DNA marker)	Query(child) DNA profile		Mother DNA profile	
	Allele_ 1	Allele_ _2	Allele_ _1	Allele_ _2
D8S1179	14	14	14	17
D21S11	28	29	28	30
D7S820	8	10	8	11
CSFIPO	10	12	10	11
D3S1358	16	16	16	19
TH01	6	9	6	9
D13S317	10	12	10	13
D16S539	11	12	11	15
D2S1338	17	23	17	18
D19S433	14	15	14	15
VWA	15	17	15	20
TPOX	8	10	8	10
D18S51	13	16	13	13
DS5818	12	12	12	15
FGA	19	23	19	27
Amelogenin	X	Y	X	X
Confidence Of matching (ARCS)				60 %

4.2 Second Scenario

The same steps of process for the previous scenario is repeated, but will be compare the confidence of matching between victim DNA profile and son DNA profile using Association Rule Classification System. Table (4.3) that shows both DNA profiles of individuals and confidence of matching value results of this scenario. The DNA profile matching result is 57% it also very acceptable. To verify the result, the comparison for this confidence of DNA profiles, set with threshold of family relationship. The similarity between an individual and his son is at least 50% according to threshold. The result produced by the proposed Association Rule Classification System is more realistic compared to if the result produced is lower than 50%. Likewise it was confirmed and verified between the experimental and real ratios, where the conformity ratio was very high. Those two profiles have a strong family relationship, so the confidence of matching result should be high.

Table (4.3) Second Scenario

Loci Name (DNA marker)	Query(victim) DNA profile		Son DNA profile	
	Allele_1	Allele_2	Allele_1	Allele_2
D8S1179	12	13	10	13
D21S11	30	33.2	29	33.2
D7S820	9	10	7	10
CSF1PO	11	12	12	12
D3S1358	16	19	18	19
TH01	6	9	68	9
D13S317	8	10	10	10
D16S539	11	12	8	12
D2S1338	19	20	18	20
D19S433	14	14	12	14
VWA	18	19	18	19
TPOX	8	9	9	9
D18S51	16	17	16	17
DS5818	12	13	10	13
FGA	22	23.2	20	23.2
Amelogenin	X	Y	X	Y
Confidence Of matching (ARCS)				57 %

4.3 Third Scenario

The aim of such scenario is to know the accuracy and validate the confidence of matching of system by comparing between two DNA profiles which do not have any family relationship. Table (4.4) that shows both DNA profiles of individuals and confidence of matching value between those profiles using Association Rule Classification System. The confidence of matching result for DNA profile is 17%. This result is lower than threshold of family relationship. This reveals that those two profiles do not have any family relationship of the two individuals referred to above, so the DNA confidence of matching result tends to be low. Likewise it was confirmed and verified between the experimental and real ratios, where the conformity ratio and accuracy was very high. Likewise, the victim's DNA profile was compared with 400 other DNA profiles that do not have any relationship with them, so the results came with a very low level of confidence of matching

for all profiles, as they are compatible with the threshold of biological family relationship.

Table (4.4) third Scenario

Loci Name (DNA marker)	Query(victim) DNA profile		Reference DNA profile	
	Allele_1	Allele_2	Allele_1	Allele_2
D8S1179	14	14	11	11
D21S11	28	29	27	30
D7S820	8	10	12	12
CSF1PO	10	12	10	11
D3S1358	16	16	16	17
TH01	6	9	7	8
D13S317	10	12	8	12
D16S539	11	12	9	13
D2S1338	17	23	18	21
D19S433	14	15	13	16
VWA	15	17	15	17
TPOX	8	10	9	11
D18S51	13	16	14	17
DS5818	12	12	11	13
FGA	19	23	25	21
Amelogenin	X	Y	X	X
Confidence Of matching (FARCS)				17 %

4.4 Fourth Scenario

The purpose of this scenario is to compare confidence of matching results between the victim DNA profiles with itself within 1000 DNA profile, with a view to know the accuracy and validate of this system and its ability to give true and realistic results. Let's see how accurate the system is in identifying the victim's profile from all other profiles. Figure (4.2) that show the confidence of matching values between those all reference and victim's DNA profiles by using Association Rule Classification System.

ID_Source	Sample Source name	ID_Refrence	Sample Reference name	Percentage Of Matching	Gender Of Matching	NO. loci matched
2	N1	1	mother 1	100 %	NOT	15
2	N1	2	son 2	17 %	NOT	5
2	N1	3	father 3	23 %	NOT	7
2	N1	4	son 4	27 %	NOT	8
2	N1	5	father 5	20 %	NOT	6
2	N1	6	son 6	33 %	NOT	10
2	N1	7	father 7	33 %	NOT	10
2	N1	8	son 8	33 %	NOT	9
2	N1	9	father 9	10 %	NOT	3
2	N1	10	son 10	33 %	NOT	9
2	N1	11	father 11	33 %	NOT	8
2	N1	12	mother 12	23 %	NOT	7
2	N1	13	daughter 13	27 %	NOT	7
2	N1	14	daughter 14	20 %	NOT	6
2	N1	15	mother 15	27 %	NOT	8
2	N1	16	mother 16	40 %	NOT	11
2	N1	17	father 17	20 %	NOT	6

Figure (4.2) Fourth Scenario dataset

The confidence of matching result for victim DNA profile is 100% that match with one only from other 1000 DNA profile. This reveals it the results that produced by the proposed Association Rule Classification System is more realistic and

will be judge that the individual who received a 100% matching ratio is the same victim person. Figure (4.4) that show the evaluation of the proposed system for all scenarios.

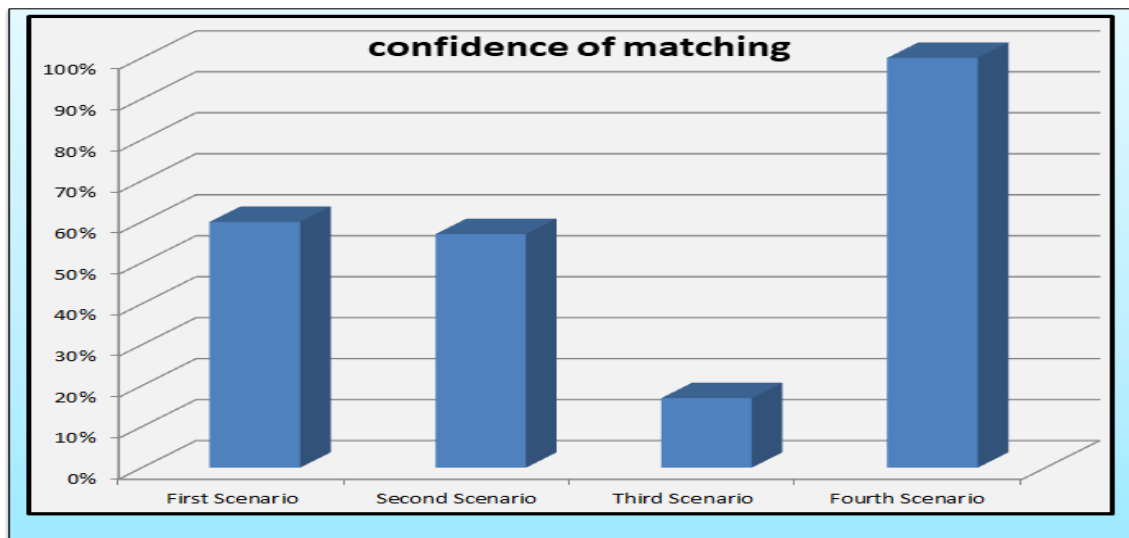


Figure (4.3) Evaluation of the Proposed System of all scenario.

5. Conclusions

The proposed biometrics system has come out with promising results in terms of human identification, database management. This is to offer a simple, automatic, reliable, and accurate forensic human identification system. ARCS using STR DNA marker for forensic identification purpose was proposed. This system provided satisfactory results in all aspects in terms of time, efficiency, capacity and accuracy as well as it adopted a STR marker, reliable and valuable tool for human recognition purpose. The introduced systems were

exploited to produce identification system based on Association Rule Classification Algorithm as Identification and biometrics technologies. This provided the proposed system with high flexibility and authenticity due to the cross checking of collected samples from crime scenes. Moreover, Visual Studio 2015, and SQL Server program was applied to build that database. The environment was used to build graphical user interface models that allow the average user to manipulate the DNA database without the need for prior knowledge of the mechanism of storing and arranging data.

REFERENCES

- [1] Faghri, Faraz, et al. "Toward scalable machine learning and data mining: the bioinformatics case." *arXiv preprint arXiv:1710.00112* (2017).
- [2] Roiger, Richard J. Data mining: a tutorial-based primer. CRC press, 2017.
- [3] Khoo, Yik-Herng, et al. "Multimodal biometrics system using feature-level fusion of iris and fingerprint." Proceedings of the 2nd International Conference on Advances in Image Processing. 2018.
- [4] Shkedy, Gary, and Dalia Shkedy. "System for Adaptive Teaching Using Biometrics." U.S. Patent Application No. 15/988,114.2019 .
- [5] Kulkarni, Sumit, and Manali Pandit. "Biometric recognition system based on dorsal hand veins." *Int J Innov Res Sci Eng Technol* 5.9 (2016): 18899-18905.
- [6] Lister, Maia. "Lacking Regulated Policy for DNA Evidence." *Themis: Research Journal of Justice Studies and Forensic Science* 6.1 (2018): 14.

- [7] Goodwin, William, Adrian Linacre, and Sibte Hadi. An introduction to forensic genetics. Vol. 2. John Wiley & Sons, 2011
- [8] Butler, John M. Forensic DNA typing: biology, technology, and genetics of STR markers. Elsevier, 2005.
- [9] Goodwin. *Forensic DNA typing protocols*. Ed. William Goodwin. Humana Press, 2016.
- [10] Shrivastava, Pankaj, Toshi Jain, and Veena Ben Trivedi. "DNA fingerprinting: a substantial and imperative aid to forensic investigation." *Eur J Forensic Sci* 3.3 (2016): 23-30.
- [11] Elsie, B. Hebsibah, et al. "Comparison of DNA from different oral swabs and its application in DNA profiling." *World J Pharm Pharm Sci* 6.5 (2017): 791-803.
- [12] Douglas, Mandy, et al. "An overview of steganography techniques applied to the protection of biometric data." *Multimedia Tools and Applications* 77.13 (2018): 17333-17373.
- [13] Widyanto, M. Rahmat, Reggio N. Hartono, and Nurtami Soedarsono. "A novel human STR similarity method using cascade statistical fuzzy rules with tribal information inference." *International Journal of Electrical and Computer Engineering* 6.6 (2016): 3103.
- [14] Croock, M. S., and S. Dh Khudhur. "Dental X-Ray Based Human Identification System for Forensic." *Engineering and Technology Journal* 35.1 Part (A) Engineering (2017): 49-60.
- [15] Anggreainy, Maria Susan, et al. "Family relation and STR-DNA matching using fuzzy inference." *International Journal of Electrical & Computer Engineering* (2088-8708) 9.2 (2019).
- [16] Han, Jiawei, Micheline Kamber, and Jian Pei. "Data mining concepts and techniques third edition." *Morgan Kaufmann* (2011).