



# Identification of a Genuine User/Imposter by Keystroke Dynamics Dataset using Machine Learning Techniques

R.Abinaya<sup>1</sup>, AN.Sigappi<sup>2</sup>

Research Scholar, Department of Computer Science and Engineering, Annamalai Nagar, India<sup>1</sup>

Associate Professor, Department of Computer Science and Engineering, Annamalai Nagar, India<sup>2</sup>

Abinayamalar@gmail.com<sup>1</sup>, aucse\_sigappi@yahoo.com<sup>2</sup>

**Abstract:** Keystroke dynamics has been used to strengthen password-based user authentication systems by considering the typing characteristics of legitimate users. Dependence on computers to store and process sensitive information has made it necessary to secure them from intruders a behavioral biometric, keystroke dynamics flow which makes utilization of the typing style of an individual can be utilized to reinforce existing security systems adequately and inexpensively and the examination of keystroke validation, to use the Discrete Cosine Transform (DCT) to describe the keystroke progression, and gives BeiHang keystroke dataset comes about are one of the well-known classifier random forest classifier it best results achieved were respectively 90% accuracy when compared with other classifier results such as Support Vector Machine and Random tree classifier.

**Keywords:** Keystroke Dynamics, Biometrics, Discrete Cosine Transform, Support Vector Machine, Random forest classifier, Random Tree Classifier, Receiver Operating Characteristic

## I. INTRODUCTION

Increase in the number of software and devices for hacking and cracking causes gains in unauthorized access which results in manipulation of important data. Methods like user ID and password which is mostly used as security is now not reliable and secure due to rapidly increase in hackers and crackers. Also, this method no longer provides consistent security measures because passwords are prone to shoulder surfing and passwords can also be hacked. To gain secure and efficient access either user must change his password frequently or the user should use the strong password (combination of alphabets, numeric and special symbols). Users do not respect these conditions as they feel them quite strict and difficult to be applied. The solution to above said problems is keystroke dynamics. Keystroke Dynamics is a behavioral biometric approach to enhance the computer access rights. It verifies the individual by its keystroke typing pattern. Keystroke biometric depends on the supposition that the composing example of every client is interesting. The target of this survey paper is to condense

the outstanding methodologies utilized as a part of keystroke flow [1].

Approaches of User Authentication: *Object Based, Knowledge Based, and Biometric Based*: Two categories are: *Physiological biometrics* and *Behavioural biometrics*. *Physiological Biometrics*: It illustrates those features that describe who the user is depending on the physical attributes e.g. fingerprints, Iris and retina scanning. For this additional hardware required. *Behavioural Biometrics*: It is based on typing pattern, Voice recognition and Signature style. Behavioral characteristics can be composed without the requirement of any extra hardware [3]. This study will focus on Behavioral biometric technique i.e. Keystroke Dynamics.

One of the recent Biometrics Technology used in upcoming research is *Keystroke Dynamics* This method analyzes the way a user types on a fatal, by monitoring the keyboard input. Since the input device is the remaining Keyboard, this approach is not exclusive. Outline of Keystroke Dynamics This technique concentrates on the composing example of a client at a lethal and after that assessing the info



distinguishing ongoing typing beat design. Keystroke dynamics are typically gotten utilizing the planning particulars of the key down or key hold or occasions. It is referred to by various names, for example, writing biometrics and composing rhythms. The principle favourable position of utilizing keystroke dynamics is that it doesn't require any additional equipment [4] Two basic features used for keystroke dynamics are *Key Hold time* and *Inter Key time*

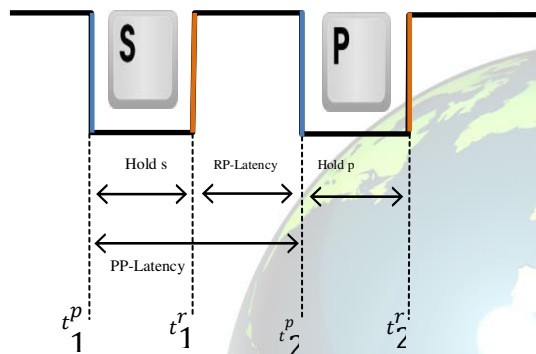


Fig 1. KEYSTROKE DYNAMICS EVENTS (PRESS AND RELEASE) AND HOLD TIME, RELEASE-PRESS (RP) LATENCY AND PRESS-PRESS (PP) LATENCY.

## II. RELATED WORKS

From the time the idea of Keystroke Dynamics was presented, much progression in the field has occurred. Several techniques came into existence then. They are described below in details with their strengths and limitations as follows: N. Chourasia Nandini, 2014 identified an additional layer of security for the authentication of the user, Keystroke Dynamics.

The security can be implemented in android phones or any other smart phones through which internet is accessible and additionally, online exchanges can be performed. Data was gathered to quantify the execution and assessment strategy was produced. A scientific model was displayed before execution. A. K. Hussain and M. M. Alnabhan, 2014 in his examination introduced a propelled keystroke authentication model increase the strength.

The keystroke structure included two segments, Firstly the deviation in typing time of client Secondly a novel client mystery code. This framework tackled the issue of huge deviations in keystroke elements and enhanced keystroke validation level was given [6]. K. Senathipathi, Krishnan

Batri, 2014 A near investigations of Particle Swarm Optimization and Genetic calculation has been demonstrated by the writer regarding keystroke dynamics.

The author select the feature selection for the proposed strategy utilizes Particle Swarm Optimization (PSO) calculation and Typing rhythms are the rawest type of information originating from the cooperation amongst clients and PCs at that point examined and broke down, they may turn into a valuable device to learn individual personality. this paper, influenced a similar examination of Particle To swarm Optimization and Genetic Algorithm as for Keystroke Dynamics [7].

The writer T. Maheswari and S. Anitha, 2014 has presented a novel approach for verification that depended on biometric attributes i.e. Keystrokes of the secret password selection. The author has measured three stages, in particular, fingerprint, login certification in view of username and password and keystroke dynamics. Two phases were likewise viewed as that are Training and testing stage. Preparing stage was actualized during enrolment and testing during verification stage [8].

D. Rudrapal, S. Das, and S. Debbarma, 2014 has combined different modules networks and figuring was performed to discover keystroke latency the as the measure of disorder. The author makes authentication and verification powerfully more secure than the typical password utilized as a part of both offline and online transactions with the assistance of experimental information. This works creator takes the Keystroke latency and duration is inadequate for user authentication, which motivates exploring other matrices. Combinations of various frameworks and figuring of level of disorder on keystroke latency and additionally term to produce client profile. Factual investigation on these lattices assesses enhanced validation process respectively. The consequence of proposed technique indicated FRR of 8% and FAR of 2%, which improved the current confirmation result utilizing keystroke dynamics [9]

A. Ahmed and I. Traore, 2014 introduced another approach for the free content investigation of keystrokes that combined monograph and digraph examination. A neural system had been utilized to anticipate missing digraphs in view of the connection between the observed keystrokes. The heterogeneous test included 53 clients, the subsequent test in a homogeneous domain considered just 17 volunteers. The results gotten from this scientist were promising with decreased error rates [10]



### III. PROPOSED RESEARCH WORK

In this project, the authentication of person by keystroke dynamics by BeiHang keystroke dynamics Datasets they proposed DCT feature vector, random forest and random tree, SVM as a classifier

#### A. The BeiHang Keystroke Dynamics Database

It can be utilized by scientists to test their algorithms and in the long run support the improvement of keystroke dynamics. There are 209 subjects engaged with building the databases. It ought to be noticed that 10 subjects of Dataset A of Database 2 are from Dataset B of Database 2. The principal database, named BeiHang Keystroke Dynamics Database 1, is caught by the online framework, and the second one, named BeiHang Keystroke Dynamics Database 2, is gathered from the implanted framework. The subjects assemble enlistment information from genuine clients utilized as preparing tests, log-in information from certified clients and log-in information from intruders. Every data are stored in text format ; they can be downloaded at. In every organizer of Database 1, the training record contains 4 or 5 enrolment tests and the file name is in the format of, say, [12345] (-regliao Xiaoying).txt meaning this is the training file for ID being 12345 and password being [liao Xiaoyin] with being the label of the file.

All the testing files have the same format: [Year-Month-Day Hour.Min.Sec] ID(-loginPSW)\_IsGenuine\_IsPostive.txt, where IsGenuine = 0 or 1 represents the data from a genuine user or an intruders; IsPostive = y or n speaks to the positive information from a client or the negative information from an intruders. For instance, the testing file, [2009-12-30 14.07.01]12345(-loginliao Xiaoying)\_1\_n.txt, demonstrates that the login time is 2009-12-30 14.07.01, ID is 12345, PSW is liao Xiaoying, and it is negative data from an intruders. The file names in Database 2 have been simplified. The data folders are named as PSW or the time when the information database were gathered. In the folder,[.txt] stores genuine client enrollment data. The whole testing records are as time-index\_IsGenuine\_ IsPostive.txt. The BeiHang Keystroke Dynamics Database 1 incorporates 1902 test tests and 477 preparing training samples from 117 subjects.

The entire Database 1 is isolated into two subsets, Dataset A and Dataset B, gathered from two unique situations. Dataset A was gathered in Internet Cafe. It contains 49 subjects, 212 training samples, 157 testing tests from bona fide clients and 996 testing tests from interlopers, as appeared in Table 1. The created business framework was inserted into the login arrangement of an online application. In Database 1, Dataset B was collected online in a university lab.

It contains 68 subjects, 265 training samples, 214 testing samples from genuine users and 535 testing samples from intruders. The BeiHang Keystroke Dynamics Database 2 was collected by the embedded system, which contains 5089 test samples and 478 training samples from 92 subjects. Dataset A and B in Database 2 are released for research purpose. Dataset A of Database 2 contains 52 subjects, 228 training samples, 717 testing samples from genuine users and 1468 testing samples from intruders. Dataset B of database 2 contains 50 subjects, 250 training samples, 1103 testing samples from genuine users and 1801 testing samples from intruders. The details are given in Table 1. It is important that there are 10 subjects seem both in these two subsets. Which contain data of stable typing rhythms? Table 1: Description of dataset A and dataset B DATASET A DATASET B Number of inducers 816 Number of inducers 365 Number of Training 417 Number of training 685 Number of Users 198 Number of users 551 Total 1428 total 1601 All the data in these dataset are original collected , without any manual modification .Generally a password is represented by following stream P1,R1,P2,R2,...,Pn, Rn, where P1 and R1 represented as the press and release time of the Ith keystroke of a password the importance of various files are appeared by their file names.

TABLE.1 DESCRIPTION OF DATASET A AND DATASET B

DATASET A		DATASET B	
Number of inducers	816	Number of inducers	365
Number of Training	417	Number of Training	685
Number of Users	198	Number of Users	551





Total	1428	Total	1601
-------	------	-------	------

Suppose a password is represented by the following sequence

$$P_1, R_1, P, R_2, \dots, P_n, R_n \quad (1)$$

All the data in these dataset are original collected , without any manual modification .Generally a password is represented by following stream  $P_1, R_1, P_2, R_2, \dots, P_n, R_n$ , where  $p_1$  and  $R_1$  represented as the press and release time of the  $i$  th keystroke of a password the meaning of different files are shown by their file names.

#### B. Database access

To download the databases for research purpose, one can visit <http://www.vmonaco.com/keystrokes> The BeiHang keystroke Dynamics -datasets or send an email to the corresponding author.

#### C. Benchmark algorithms

The framework of our Keystroke Dynamic System is shown in Fig.2 Feature extraction and classification algorithm are the main components and are discussed in detail in the following sections.

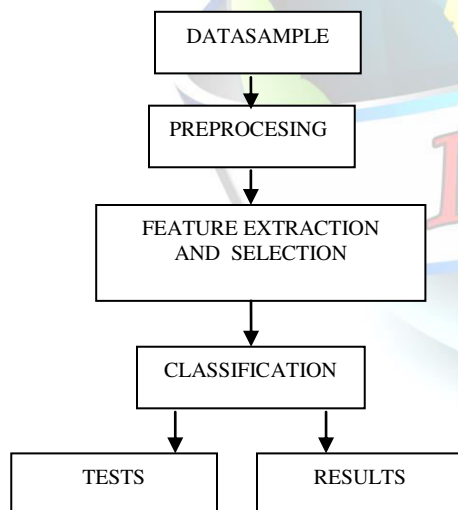


Fig. 2 FRAMEWORK OF THE KEYSTROKE DYNAMICS SYSTEM

#### D. Feature Extraction

Where represent the press and the release time of the  $i$ th keystroke of a password. The elements of the feature vector extracted from the original keystroke information are classified into two categories: dwelling time and flight time. The dwelling time is calculated by  $R_i - P_i$ , and the flight time by  $P_i - R_{i-1}$ .

Therefore, the extracted feature from the original sequence is represented as:

$$I = (R_1 - P_1, P_2 - R_1, R_2 - P_2, \dots, P_n - R_{n-1}, R_n - P_n). \quad (2)$$

The above feature is also called the original feature. The number of the registration samples collected in the training procedures.

#### E. Discrete Cosine Transforms(DCT):

In signal and image processing, DCT is extensively used to transform the whole image for feature extraction by separating the relevant coefficients and performs energy compaction .The DCT consists of three components frequencies i.e. low, middle, high each contains more significant information of an image. The low frequency usually contains the mean intensity of an image which is the most projected in FR systems [11][12]. Mathematically, the 2D-DCT of an image is given by

$$F(u, v) = \alpha(u)\alpha(v) \sum_{x=0}^{N-1} \sum_{y=0}^{M-1} \cos \left[ \frac{\pi u}{2N} (2x+1) \right] \cos \left[ \frac{\pi v}{2M} (2y+1) \right] f(x, y) \quad (3)$$

$$\alpha(u)\alpha(v) = \begin{cases} \sqrt{\frac{1}{N}} & \text{for } u, v \neq 0 \\ \sqrt{\frac{2}{N}} & \text{for } u, v = 0 \end{cases} \quad (4)$$

Where  $M \times N$  is the size of the image, where  $f(x, y)$  is the intensity of the pixel at coordinates  $(x, y)$ ,  $u$  varies



from 0 to  $M-1$ , and  $v$  varies from 0 to  $N-1$ , This project DCT used to varying the Dimension value

#### F. Feature vector

By these sequence the password is represented as

$$P_1, R_1, P_2, R_2, \dots, P_n, R_n$$

Where  $P_i$  is press time and  $R_i$  is release time of the  $i$ th keystroke of a password and the feature vector extracted from the original keystroke information are classified in to two categorizes: Dwelling time and Flight time. Dwelling time is calculated by  $P_i - P_{i-1}$  and the Flight time is calculated by  $R_i - P_i$

#### G. Feature Vector Extraction:

Therefore the extracted feature from the original sequence is represented as

$$I = (P_1 - R_1, R_1 - P_2, \dots, P_n - R_{n-1}) \quad (5)$$

By using this sequence. We get Feature Vector  $V = 117340, 165499, 71146, 205008, 88286, 357374, 91145, 308217, 85432, 22938, 102572$

Supposed we have  $k$  classifiers, whose classifiers and DCT feature combination score is denoted as score level fusion vector.

#### H. Classification Algorithms:

##### Fusion of Features and Classifiers:

Mixture of different knowledgeable decision of widely examined in past twenty years. Combination strategies can be gathered by the level at which they work. The least complex route is in the feature level, where various types of highlights are linked into an expanded feature vector. This combination acquires the benefits of various highlights, and any classifier is effectively utilized with them to fabricate the last classification display. Combination should likewise be possible in the level of choice or output score, which is called classifier-level combination. It is a very prominent path as the score is for the most part considered as another sort of feature. This paper researches the two techniques for execution performance metrics. For include level

combination, we can without much of a stretch get the new expanded feature. DCT feature vector Similar to the component level combination, the classifier-level blend depends on the scores of classifiers.

Supposed it have  $k$  classifiers, whose classifiers and DCT feature combination score is denoted as score level fusion vector

#### I. Experiment:

In this area, we exhibit benchmark test comes about for some grouping and feature extraction calculations on the BeiHang Keystroke Dynamics Databases. The broad research tests bring the assessment of various features, classifiers, and their combinations.

We also exhibit those specific rhythms for different individuals can lead to high performance, which can be used in practical applications, such as password protection

#### J. Evolution criteria:

In the research experiments, we utilize the False Positive Rate and the True Positive Rate for assessment measurements. The previous is the level of intruders who can enter the account by emulating the typing style of genuine clients.

The last is the level of genuine clients who can effectively sign into the framework with the correct keystroke way. By changing the edge in the grouping methodology, we get a progression of True Positive Rates and False Positive Rates, and then we use these results to draw a ROC curve. The ROC curve is used for evaluation of various algorithms including the Random forest classifier, classifier random tree classifier with the original feature DCT. We also provide the Equal Error Rate (EER) for further evaluation of different methods. EER is the percentage where the False Positive Rate equals the False Negative Rate

#### K. Classifier, Features, Training:

This section explains the classifier that we used, the features it employed, and its training and testing. The MATLAB Programming environment (version 2013 a) was used for analyses And WEKA 6 for Classifier analysis.



#### *L. Classifier – Random Forest classifiers:*

Random Forest is one of the most versatile machine learning algorithms available in the network. With this inherent assembling limit, the errand of building a conventional summed up display (on any dataset) gets substantially less demanding. Be that as it may, I've seen individuals utilizing random forest as a (black box) discovery demonstrates i.e., they don't comprehend what's going on underneath the code. They simply code. Truth be told, the simplest piece of machine learning is coding. In the event that you are new to machine taking in, the random forest algorithm calculation ought to be on your tips. Its capacity to take care of both regression and classification issues alongside strength to related features and variable significance plot gives us enough make a beeline for take care of different issues. Each tree is developed as takes after

1. On the off chance that the quantity of cases in the preparation set is  $N$ , test  $N$  cases indiscriminately - yet with substitution, from the original database information. This specimen will be the preparation set for developing the tree.

2. On the off chance that there are  $M$  input factors, a number  $m \ll M$  is indicated with the end goal that at every hub,  $m$  factors are chosen aimlessly out of the  $M$  and the best split on these  $m$  is utilized to part the hub. The estimation of  $m$  is held steady amid the random forest developing.
3. Each tree is developed to the biggest degree conceivable. There is no pruning. In the first paper on arbitrary random forest classifiers, it was demonstrated that the randomly mistake rate relies upon two things: The relationship between's any two trees in the backwoods. Expanding the relationship builds the random forest mistake rate. The quality of every individual tree in the random forest. A tree with a low blunder rate is a solid classifier. Expanding the quality of the individual trees diminishes the random forest mistake rate. Diminishing  $m$  decreases both the connection and the tests. Expanding it increments both. Some place in the middle of is an "ideal" scope of  $m$  - as a rule very wide. Utilizing the error rate (see underneath) an estimation of  $m$  in the range can rapidly be found. This is the main movable parameter to which irregular woods is to some degree delicate.

#### *M. Features used in the classifiers:*

Features utilized as a part of the classifier during typing, all key-press (key-down) and key-release (key-up) events were time stamped and recorded.

From these events, each of the three features used in the random forest classifier, SVM, Random tree can be derived: (1) hold time (time elapsed from key-down to key-up of a single key); (2) digram latency (time elapsed from the key-down of a character being typed to the key-down of the next character); and (3) diagram interval (key-up to key-down latencies between diagrams). For a ten-digit passcode, there are 11 hold times (including the return key), 10 key-down to key-down latencies, and 10 key-up to key-down intervals, which taken together form a 31-dimensional vector that represents each passcode repetition. All classifier features were used; because they form a superset of the features commonly used by other researchers. It is unexcelled in precision among current calculations. It runs effectively on extensive databases. It can deal with a large number of info factors without variable detection. It gives appraisals of what factors are essential in the classification. It produces an inner fair-minded gauge of the speculation error as the forest building advances. It has a viable technique for assessing missing information and keeps up exactness when an extensive extent of the information is absent. It has strategies for adjusting blunder in class populace lopsided informational indexes. The created random forest can be put something aside for later use of other information.

Models are registered that give data about the connection between the factors and the classification. It registers vicinities between sets of cases that can be utilized as parts of bunching, finding anomalies or (by scaling) give intriguing perspectives of the information. The capacities of the above can be stretched out to unlabeled information, prompting unsupervised grouping, information perspectives, and exception location. It offers an experimental method for detecting variable interactions. Although some of these features are linearly dependent, this is not a concern when using a random forest, because the random forest performs feature selection as part of its training, thereby accommodating any linear dependencies among features





#### IV. RESULT AND DISCUSSION:

##### DATASET A:

##### Results for dataset A

TABLE. 2 ACCURACY AND EER RESULTS WITH DIFFERENT FEATURE DIMENSIONS AND CLASSIFIER ON DATASET 1, WHERE\* INDICATES RESULTS

Dataset A	Random Forest Classifier		Random tree Classifier		Support Vector Machine classifier	
	Accuracy	EER Rate	Accuracy	EER Rate	Accuracy	EER Rate
DCT – Feature (DIM=10)	88.7	11.3	81.0	18.9	79.6	20.4
DCT – Feature (DIM=12)	89.4	10.6	80.5	19.4	83.2	16.8
DCT – Feature (DIM = 14)	85.9	14.1	80.1	19.8	86.3	13.7

##### DATASET B

##### Results for dataset B:

TABLE.3 : ACCURACY AND EER RESULTS WITH DIFFERENT FEATURE DIMENSIONS AND CLASSIFIER ON DATASET B, WHERE\* INDICATES RESULTS

Dataset A	Random Forest Classifier		Random tree Classifier		Support Vector Machine classifier	
	Accuracy	EER Rate	Accuracy	EER Rate	Accuracy	EER Rate
DCT – Feature (DIM=10)	87.5	12.5	79.4	20.5	81.2	18.3

DCT – Feature (DIM=12)	88.8	11.2	80.5	19.4	78.8	21.2
DCT – Feature (DIM = 14)	82.5	17.5	81.6	18.3	84.3	15.7

##### ROC CURVE:

Receiver operating curve or Relative operating trademark (ROC): The ROC plot is a visual portrayal of the exchange off between the FMR and the FNMR.

When all is said and done, the coordinating calculation plays out a choice in light of a limit that decides how near a format the information should be for it to be viewed as a match. This kind of chart is known as a Receiver Operating Characteristic bend (or ROC bend.) It is a plot of the genuine positive rate against the false positive rate for the diverse conceivable cut purposes of an indicative test. A ROC bend exhibits a few things: It demonstrates the tradeoffs amongst affectability and specificity (any expansion in affectability will be joined by a reduction in specificity). The nearer the bend takes after the left-hand fringe and afterward the best outskirt of the ROC space, the more precise the test. The nearer the bend goes to the 45-degree corner to corner of the ROC space, the less precise the test.

The incline of the digression line at a cut point gives the probability proportion (LR) for that estimation of the test. You can look at this on the chart above. Review that the LR for  $T4 < 5$  is 52. This compares to the far left, soak bit of the bend. The LR for  $T4 > 9$  is 0.2. This relates to the far right, about flat part of the bend. The region under the bend is a measure of content precision. Comparison of ROC results for two datasets A and B

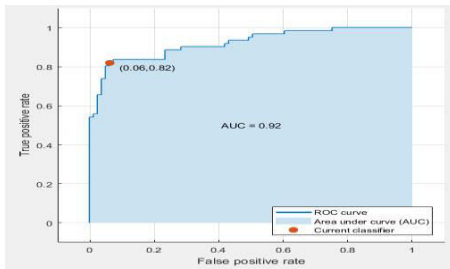


Fig. 3 Dataset A in Dimension 10\_Random forest Classifier  
\_ROC\_Accuracy \_88.7

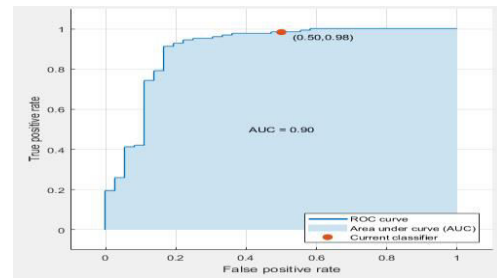


Fig. 6 Dataset B in Dimension 10\_Random forest Classifier  
\_ROC\_Accuracy \_87.5

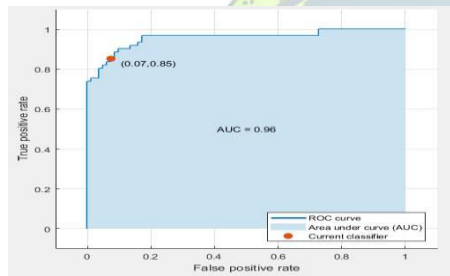


Fig. 4 Dataset A in Dimension 12\_Random forest Classifier  
\_ROC\_Accuracy \_89.4

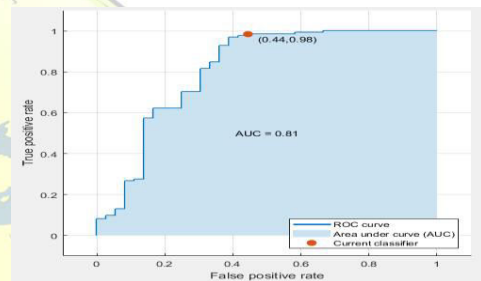


Fig.7 Dataset B in Dimension 12\_Random forest Classifier  
\_ROC\_Accuracy \_88.8

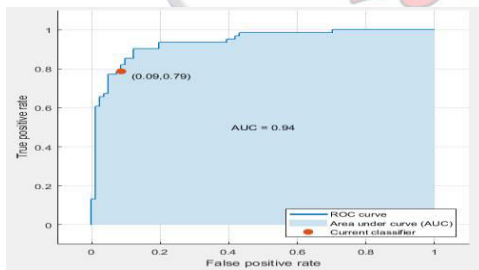


Fig.5 Dataset A in Dimension 14\_Random forest Classifier  
\_ROC\_Accuracy \_85.9

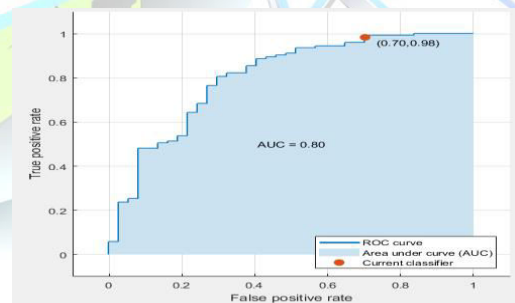


Fig .8 Dataset B in Dimension 14\_Random forest Classifier  
\_ROC\_Accuracy \_82.5





#### Performance Evolution:

The execution of a biometric framework is by and large described by the receiver operating characteristic (ROC). It can be condensed by the equal error rate (EER) the point on the bend where the false acceptance rate (FAR) and false rejection rate (FRR) are equivalents. Other framework assessment criteria incorporate productivity, flexibility, convenience, and comfort. Performance measures the execution of a biometric framework as far as procurement and recognizing error.

The end goal to assess the execution of a biometric framework, we by and large need a test benchmark and execution measurements. As per the International Organization for Standardization ISO/IEC 19795-1, the execution measurements are partitioned into three sets: Acquisition execution measurements, for example, the Failure-To-Enrol rate (FTE). Check framework execution measurements, for example, the Equal Error Rate (EER). Identification framework execution measurements, for example, the False-Negative and the False-Positive Identification Rates (FNIR and FPIR, separately).

Effectiveness: Effectiveness shows the capacity of a technique to accurately separate genuine and imposter. Execution pointers utilized by the inquire about are compressed as take after.

False Rejection Rate (FRR) alludes to the rate proportion between erroneously denied honest to genuine clients against the aggregate number of authentic clients getting to the system. Once in a while known as False Nonmatch Rate (FNMR) or sort 1 error. A lower FRR infers less rejection rate and less access by genuine users. False Acceptance Rate (FAR) is characterized as the rate proportion between dishonestly acknowledged unapproved clients against the aggregate number of intruders getting to the framework. Terms, for example, False Match Rate (FMR) or sort 2 blunder alludes to a similar significance. A littler FAR shows an imposter accepted. Equal Error Rate (EER) is utilized to decide the general exact accuracy and in addition a similar estimation against different frameworks. It might be here and there referred for as Crossover Error Rate (CER). Result examination depicted in the following segment will basically be express with FAR, FRR, and EER

$FRR = (\text{Number of refused genuine}) / (\text{Total no of genuines})$

$FAR = (\text{Number of Accepted imposter}) / (\text{Total no of imposter})$

TABLE .4 A10 CONFUSION MATRIX– A10 - RANDOM FOREST (ACCURACY = 89.3)

Classification		
Actual class	516	96
	57	759

TABLE .5 A12 – CONFUSION MATRIX A12 - RANDOM FOREST (ACCURACY = 90.5)

Classification		
Actual class	525	87
	48	768

TABLE .6 A14 – CONFUSION MATRIX A14 - RANDOM FOREST (ACCURACY = 85.5)

Classification		
Actual class	467	142
	64	752

Table .7 B10 – Confusion Matrix B10 - RANDOM FOREST (accuracy = 87.8)

Classification		
Actual class	490	132
	44	772



TABLE .8 B12 – CONFUSION MATRIX B12 - RANDOM FOREST  
(ACCURACY = 88.1)

Classification		
Actual class	512	100
	74	742

TABLE .9 B14 – CONFUSION MATRIX B14 - RANDOM FOREST  
(ACCURACY = 82.4)

Classification		
Actual class	426	186
	65	751

Where P represents positive rate and N represents Negative rate, Where Classification means a correct classification of the instances. Originating from a deceptive messages are instances that are supposedly truthful but classified as deceptive are instances that were derived from deceptive messages but are classified as truthful. are instances that originate from truth messages and are classified as such some of actual deceptive (positive) instances is p with  $p = \frac{512}{512+74}$  here random forest classifiers takes B14 dimension get  $p = 88.1$  and the sum of actual truthful (negative) instances N with  $N = \frac{100}{100+742}$  and P and N with  $p = \frac{512}{512+100}$ ,  $p = 83.2$  N =  $\frac{100}{100+742}$ , N=855 represented as sum of instances are classified respectively

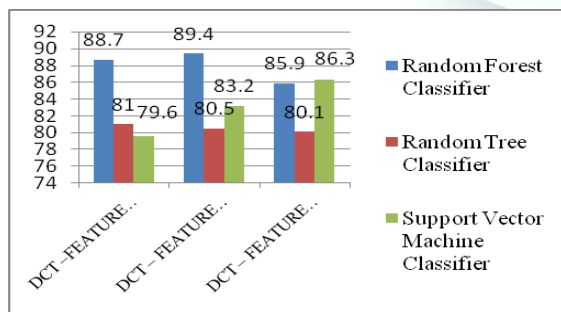


Fig. 9 Different Dimensions Level by DCT for various classifiers Dataset A

The above chart describes the different dimension level by DCT accuracy rate with various classifiers. Hence the performance matrix shown by the ROC curve

## V. CONCLUSION

Two large databases have been collected and open for open research. Different features and benchmark calculations have been tested and outlined. We designed both DATASETS A, B for security device and an online keystroke dynamics system. The new component incorporates DCT and their combination. The benchmark comes about are gotten by the Random Forest classifier display as the classifier, applied on the original and broadened features. Our future work will concentrate on boosting the classifiers and promoting the applications. Thus the Random forest classifier got the better outcome with 90% accuracy respectively.

## REFERENCES

- [1] R.Abinaya, AN.Sigappi 2016 "Survey of Biometric Systems: Evolution and Challenges" Multimedia signal processing and Applications (NCMSPA) , ISBN No. 978-81-922221-6-5., pg.no.139.
- [2] Pin Shen Teh,1 Andrew Beng Jin Teoh and Shigang Yue1, 2013 "A Survey of Keystroke Dynamics Biometrics " Hindawi Publishing Corporation The Scientific World Journal Volume, Article ID 408280,Pg.no .24
- [3] Shimaa I. Hassan 1 , Mazen M. Selim2 , and Hala H. Zayed 2013 "User Authentication with Adaptive Keystroke Dynamics" IJCSI International Journal of Computer Science Issues, Vol. 10, Issue 4, No 2, 2013 ISSN (Print): 1694-0814
- [4] Roy A. Maxion and Kevin S. Killourhy , 2010 "Keystroke Biometrics with Number-Pad Input" IEEE/IFIP International Conference on Dependable Systems & Networks (DSN)
- [5] Rosy Vinayak1 , Komal Arora, 2015 "A Survey of User Authentication using Keystroke Dynamics " International Journal of Scientific Research Engineering & Technology (IJSRET), ISSN 2278 – 0882 Volume 4, Issue 4.
- [6] .Hussain et al. (A. K. Hussain and M. M. Alnabhan, 2014 " Advanced authentication scheme using a predefined keystroke structure" International Journal of Computer 2014
- [7] K. Senathipathi, Krishnan Batri, , 2014 "An analysis of Particle Swarm Optimization and Genetic Algorithm with respect to keystroke dynamics" Green Computing Communication and Electrical Engineering (ICGCCEE), Electronic ISBN: 978-1-4799-4982-3 ,DVD ISBN: 978-1-4799-4983-0
- [8] Senathipathi et al.(K. Senathipathi, Krishnan Batri, "An analysis of Particle Swarm Optimization and Genetic Algorithm with respect to keystroke dynamics" Green Computing Communication and Electrical Engineering (ICGCCEE), Electronic ISBN: 978-1-4799-4982-3 ,DVD ISBN: 978-1-4799-4983-0 , 2014
- [9] Maheswari et al. (T.Maheswari and S. Anitha, 2014) "Remote User Authentication Scheme: A Comparative Analysis and Improved Behavioral Biometrics Based Authentication Scheme"Published



- in:** Micro-Electronics and Telecommunication Engineering (ICMETE), 2016
- [10] Rudrapal et al. (D. Rudrapal, S. Das, and S. Debbarma, 2014) *"Improvisation of Biometrics Authentication and Identification through Keystrokes Pattern Analysis"* International Conference on Distributed Computing and Internet Technology ICDCIT 2014: Distributed Computing and Internet Technology pp 287-292
- [11] Ahmed et al. (A. A. Ahmed and I. Traore, 2014) *"Biometric Recognition Based on Free-Text Keystroke Dynamics"* **Published in:** IEEE Transactions on Cybernetics ( Volume: 44, Issue: 4, April 2014 ) **Referenced in:** IEEE Biometrics Compendium IEEE RFIC Virtual Journal IEEE RFID Virtual Journal, **Page(s):** 458, 472

