



HAL
open science

Soft Biometrics for Keystroke Dynamics: Profiling Individuals While Typing Passwords

Syed Zulkarnain Syed Idrus, Estelle Cherrier, Christophe Rosenberger, Patrick Bours

► **To cite this version:**

Syed Zulkarnain Syed Idrus, Estelle Cherrier, Christophe Rosenberger, Patrick Bours. Soft Biometrics for Keystroke Dynamics: Profiling Individuals While Typing Passwords. Computers & Security, 2014, pp.1. hal-01011801

HAL Id: hal-01011801

<https://hal.science/hal-01011801>

Submitted on 24 Jun 2014

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Accepted Manuscript

Soft Biometrics for Keystroke Dynamics: Profiling Individuals While Typing Passwords

Syed Zulkarnain Syed Idrus, Estelle Cherrier, Christophe Rosenberger, Patrick Bours



PII: S0167-4048(14)00089-3

DOI: [10.1016/j.cose.2014.05.008](https://doi.org/10.1016/j.cose.2014.05.008)

Reference: COSE 805

To appear in: *Computers & Security*

Received Date: 31 October 2013

Revised Date: 6 April 2014

Accepted Date: 25 May 2014

Please cite this article as: Syed Idrus SZ, Cherrier E, Rosenberger C, Bours P, Soft Biometrics for Keystroke Dynamics: Profiling Individuals While Typing Passwords, *Computers & Security* (2014), doi: 10.1016/j.cose.2014.05.008.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

Soft Biometrics for Keystroke Dynamics: Profiling Individuals While Typing Passwords

Syed Zulkarnain Syed Idrus^{1,2,3,4*}, Estelle Cherrier^{2,3,4}, Christophe Rosenberger^{2,3,4}, Patrick Bours⁵,

¹*Universiti Malaysia Perlis, 01000 Kangar, Perlis, Malaysia*

²*Université de Caen Basse-Normandie, UMR 6072 GREYC, F-14032 Caen, France*

³*ENSICAEN, UMR 6072 GREYC, F-14032 Caen, France*

⁴*CNRS, UMR 6072 GREYC, F-14032 Caen, France*

⁵*NISlab, Gjøvik University College, Gjøvik, Norway*

Abstract

This paper presents a new profiling approach of individuals based on soft biometrics for keystroke dynamics. *Soft biometric traits* are unique representation of a person, which can be in a form of physical, behavioural or biological human characteristics that differentiate between him/her into a group people (*e.g.* gender, age, height, colour, race *etc.*). *Keystroke dynamics* is a behavioural biometric modality to recognise how a person types on a keyboard. In this paper, we consider the following soft traits: the hand category (*i.e.* if the user types with one or two hands), the gender category, the age category and the handedness category. For this purpose, we collected a new database. Two cases are studied: static passwords and free text. By combining machine learning and fusion process, the results are promising.

*Corresponding author

Email addresses: `syed-zulkarnain.syed-idrus@ensicaen.fr` (Syed Zulkarnain Syed Idrus^{1,2,3,4*}), `estelle.cherrier@ensicaen.fr` (Estelle Cherrier^{2,3,4}), `christophe.rosenberger@ensicaen.fr` (Christophe Rosenberger^{2,3,4}), `patrick.bours@hig.no` (Patrick Bours⁵)

1
2
3
4
5
6
7
8
9 *Keywords:* Biometrics, keystroke dynamics, soft biometrics, pattern
10 recognition, data fusion, computer security.
11
12

14 1. Introduction

15
16
17 It is accepted that the way a person types on a keyboard contains timing
18 patterns, which can be used to label him/her and this is called *keystroke*
19 *dynamics*. Keystroke dynamics is an interesting and a low cost biomet-
20 ric modality [1, 2], indeed for example no additional device is required.
21 Keystroke dynamics belongs to the class of behavioural biometrics, in the
22 sense that the template of a user reflects an aspect of his/her behaviour.
23 Among the behavioural biometric modalities, we can mention signature dy-
24 namics analysis, gait recognition, voice recognition, or keystroke dynamics
25 [3, 4, 5, 6]. In general, the global performances of behavioural biometric
26 modalities (and especially keystroke dynamics) based authentication systems
27 are lower than the popular morphologic biometric modalities based authen-
28 tication systems (such as fingerprints, face or iris)[7, 8]. The fact that the
29 performances of keystroke dynamics are lower than other biometric modal-
30 ities can be explained by the intra-class variability of the users behaviour.
31 This intra-class variability pertaining to computer users can be accounted
32 for by a way of typing which is different when they are nervous, or angry, or
33 even sad ... [9].
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

50 One solution to cope with this variability is to study *soft biometrics*, which
51 was first introduced by Jain *et al.* in [10]. In that paper ‘*soft biometric*
52 *traits*’ are defined as “*characteristics that provide some information about*
53 *the individual, but lack the distinctiveness and permanence to sufficiently*
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9 *differentiate any two individuals*". Jain *et al.* considered gender, ethnicity
10 and height as complementary data for a usual fingerprint based biometric
11 system. Thus, soft biometrics allow a refinement of the search of the genuine
12 user in the database, resulting in a computing time reduction. For example,
13 if the capture corresponds to a male according to a soft biometrics module,
14 then, the standard biometric identification system can confine its search area
15 to male users, without considering female ones.
16
17
18
19
20
21

22 Since the work of Jain *et al.*, several other articles related to soft bio-
23 metrics can be found in the literature. In the paper [11], body weight and
24 fat measurements are considered as soft criteria to enhance a standard fin-
25 gerprint based biometric system. An overview can be found in [12] about
26 soft biometrics, under a '*Bag of Soft Biometrics*', where Dantcheva *et al.*
27 make a comparison with the pioneering work of Alphonse Bertillon, whose
28 anthropometric criteria gave rise to soft biometrics [13]. This paper proposes
29 some facial soft biometrics and also body soft biometrics, namely weight and
30 clothes colour detection. In [14], Park and Jain present how gender or ethnic-
31 ity and facial marks such as scars, moles and freckles can be used to enhance
32 face recognition. In reference [15], shape based eyebrow features are used for
33 biometric recognition and soft biometric classification. In [16], the authors
34 use soft biometrics (height and colour model of head, torso and legs) to help
35 identifying people in videos in surveillance networks. Marcialis *et al.* [17] use
36 hair colour and ethnicity as soft biometrics combined with face modality.
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Regarding keystroke dynamics, Bixler and D’Mello [18] look into the likelihood of 44 people’s behaviour, whether they stay idle, involved or bored when asked to write on a given task. Their result are between 11% and 38% higher than random guessing. In our previous study [19], the results also show that it is possible to detect users’ way of typing by using one/two hand(s) with over 90% recognition rate; gender between 65% and 90%; age between 65% and 82%; and handedness between 70% and 90% correct recognition accuracy with 110 users.

The objective of this paper is to propose an extended study of soft biometrics for keystroke dynamics from our previous study in [19] on a new biometric benchmark database called ‘GREYC-NISLAB Keystroke’ [20] that we have created. We propose in this paper a thorough evaluation of the soft biometrics system and a comparison between static passwords and free text (digraphs). Thus, the novelty (compared to our papers mentioned) is to study to what extent soft biometrics can enhance the recognition performance of keystroke based authentication systems. Furthermore, we show how the performances can be increased significantly by data fusion for passwords. As soft criteria, we propose to test if it is possible to predict if the user:

1. types with one or two hands
2. is a male or a female
3. belongs to a particular age category
4. is right- or left-handed

This paper is organised as follows. Section 2 is devoted to the description of the proposed method. In Section 3, we describe the protocol that applied and present the obtained results on the benchmark database in Section 4. Section 5 presents the conclusions and the different perspectives of this study.

2. Proposed Methodology

In general, keystroke dynamics authentication systems involve a keyboard and an application for the capture and processing of the biometric information. Users are required to type on a keyboard running a dedicated application. Each capture is stored in a database within the application in the form of keystroke or timing features for all correct and incorrect entries. These features are composed of several timing values that are extracted, which is the *pattern vector* that is used for the analysis. For each soft criterion, two steps are involved in recognition evaluation: (i) a training step, and (ii) a test step, both relying on a machine learning algorithm. Here we have chosen SVM (Support Vector Machine) [21], on account of its efficiency. As a result, we compute the accuracy rate of the prediction of each soft category by the trained SVM. A graphical representation of the overall process is illustrated in Figure 1. In order to enhance the overall recognition performance, data fusion is then applied.

Figure 1: The overall process of the proposed system.

2.1. Data Capture

Different types of features can be extracted from a user while typing on a keyboard [2]: “(i) code of the key; (ii) the type of event (press or release); and (iii) the time of the event”. All this timing information is stored in the form of raw data, which contains: (see Figure 2)

- $ppTime$ (PP): the latency of when the two buttons (keys) are pressed;
- $rrTime$ (RR): the latency of when the two buttons (keys) are released;
- $prTime$ (PR): the duration of when one button (key) is pressed and the other is released;
- $rpTime$ (RP): the latency of when one button (key) is released and the other is pressed.
- $vector$ (V): the concatenation of the previous four timing values.

Figure 2: Keystroke typing features.

Subsequently, the keystroke template V is utilised for the analysis for each soft category. For keystroke dynamics systems, we apply two approaches, namely: static passwords and free text. Concerning static passwords, we analyse all the typing features previously described. For free text, the analysis is based on digraphs, which are the time latencies between two successive keystrokes *i.e.* digraphs transition time. These typing rhythms are extracted from the users' texts typed without any specific constraint.

2.2. Data Analysis

For the data analysis, we recall that we are interested in soft biometrics criteria that can be applied to our biometric database: one or two hand(s); male or female, age < 30 or ≥ 30 years old, right- or left-handed. This subsection presents the methodology in which we followed to analyse keystroke data. Classification is performed by training, for each soft criterion (hand,

gender, age, and handedness categories), using a Support Vector Machine. We use LibSVM [22] with the Radial Basis Function (RBF) kernel [23, 24]. In order to maximise the performance, we have to determine which is the best couple for our computation. We set the following values for the parameters: $C = 128$ is the penalisation coefficient of the SVM; $\gamma = 0.125$ is the parameter of the kernel, as introduced by [23].

The computation of the SVM process is repeated for 100 iterations for each percentage of the training ratio, to produce an averaged recognition rate.

2.3. Data Fusion Process

For the data fusion, we apply two techniques based on *majority voting* and *score fusion* with binary classifications. For the sake of clarity, we take the example of gender category. There are more men than women in the database (*i.e.* 78 males; 32 females). We select data to have the same number within each category, so here, we randomly remove 46 males. We keep the same users sub-sample for each password, and we train one SVM per soft category. To avoid the influence of sample extraction, the whole process (from the extra men removal to the fusion) is repeated 100 times, with a different random draw of 32 males each time. The presented results are the average of these 100 classifications. Now, we present the chosen fusion processes.

First fusion process: *majority voting*. The predicted label (+1 or -1) is exploited in the first fusion method: the majority voting. Since there are 5 passwords, the majority is easily obtained.

Second fusion process: *score fusion*. We compute this score by using the

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

predicted label and its associated probability. This method, we obtain a score in the range $[0; 1]$, then we compute the average of 5 probabilities to decide the final class. If the average is above 0.5 then the label 1 is assigned, otherwise label 0.

Once this process has been completed, we can compute the *confusion matrix* [25] to obtain the correct recognition rate for each class. To compute the recognition rate (for gender category), we apply formula (1), where $M_correct$ and $F_correct$ are respectively the total number of correctly predicted Males and Females. A large value of r guarantees a large correct recognition rate for the considered category. Subsequently, we will be able to evaluate to what extent both fusion processes can enhance the performance.

$$r = \frac{M_correct + F_correct}{total_data} \times 100\% \quad (1)$$

3. Experimental Protocol

3.1. Static Passwords

In this section, we briefly describe the protocol that we applied. We refer the interested reader to our previous paper [20] for more details. As mentioned earlier, we created a biometric benchmark database. The database can facilitate and accelerate reproducible and comparable research. In [19], an experiment has been performed in two locations: France and Norway, and a total of 110 individuals had volunteered to participate. We used two desktop keyboards (French keyboard for users in France and Norwegian keyboard for users in Norway) *i.e.* AZERTY and QWERTY (this is not a classical

1
2
3
4
5
6
7
8
9 QWERTY keyboard, however, we do not use specific Norwegian keys), re-
10 spectively. Giot *et al.* work show that the keyboard does not influence the
11 performance [26].
12
13

14
15 During the data acquisition, some metadata such as gender, age and
16 handedness were collected. We have chosen passphrases that are well-known
17 in both countries, which are between 17 and 24 characters long including
18 spaces (see Table 1). All the participants were asked to type the 5 different
19 passphrases 20 times (10 times with one hand and 10 times with two hands).
20
21
22
23
24

25
26 Table 1: Passphrases.

27
28

Label	Description	Size
Password 1 (P_1)	leonardo dicaprio	17-char
Password 2 (P_2)	the rolling stones	18-char
Password 3 (P_3)	michael schumacher	18-char
Password 4 (P_4)	red hot chilli peppers	22-char
Password 5 (P_5)	united states of america	24-char

29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

48
49
50 We have used GREYC Keystroke software developed at GREYC Labora-
51 tory (downloadable from the following address: [http://www.ecole.ensicaen.](http://www.ecole.ensicaen.fr/~rosenber/keystroke.html)
52 [fr/~rosenber/keystroke.html](http://www.ecole.ensicaen.fr/~rosenber/keystroke.html)), to capture the biometric data. At the end
53 of the data collection, a total of 11000 data samples are in the proposed bio-
54
55
56
57

metric benchmark database. For each user, 7 out of 10 samples are used for training and testing data. The first three entries for each user are not taken into account because leeway was given to the users to allow them to train themselves for each of the given passphrases. We justify why three entries have been discarded by operational reasons: (i) noticed that we made 10 captures for each password entry because we want to avoid the volunteers from being annoyed, having to type (the same text) too many times; and (ii) by removing more than 3 captures will definitely lead to a smaller database. So it is more an operational justification than a statistical one, which could have made sense with much more data.

We define two classes C_1 and C_2 for each category as follows:

- *Hand category*: $C_1 = \text{One Hand}$: only one hand is used (right/left depending on the handedness of the user); $C_2 = \text{Two Hands}$: both hands are used.
- *Gender category*: $C_1 = \text{Male}$; $C_2 = \text{Female}$.
- *Age category*: $C_1 = < 30$ years old; $C_2 = \geq 30$ years old.
- *Handedness category*: $C_1 = \text{Right-handed}$; $C_2 = \text{Left-handed}$.

Here, for hand category, we use all the data. Whereas for the other soft biometrics information, we only use data corresponding to the usual way of typing, that is 2 hands.

To validate the proposed recognition system, we compute Confidence Intervals (CI). A CI is necessary when it is associated with the recognition rate of the soft biometric trait to reinforce the confidence in the obtained results.

It represents a measure of confidence on the estimated error rate. It is based on a re-sampling, which consists of a random draw with a replacement of new values of example from the test base. For each draw, the data are randomly selected. This is done $N=100$ times in order to calculate the CI, where we perform the computation of the recognition rate for each of the N tries. The CI can be determined based on the percentiles of the normal distribution. Here, the CI at 95% is defined by Equation 2, where $m(rate)$ is the mean of the recognition rates over N iterations, and $\sigma(rate)$ is the corresponding standard deviation.

$$CI = m(rate) \pm 1.96 \frac{\sigma(rate)}{\sqrt{N}} \quad (2)$$

3.2. Free Text

Subsequently, we perform a distance measure to consider the different timing information between two-character sequences known as *digraphs*. The types passwords are considered as a whole, and only digraphs are kept. Digraphs are the latency times between two successive keystrokes. We extract the keystroke features using the mean and variance of digraphs. Here, we consider free text as the collection of the 5 passwords. Therefore, the digraphs appear with an occurrence between one and four. To obtain significative results, we restrict to digraphs with an occurrence equal or larger than 2. Thus, we consider three categories of digraph: (i) 11 with two occurrences; (ii) 2 with three occurrences; and (iii) 1 with four occurrences. Consequently, there are a total of 14 occurrences as shown in Figure 3.

In some instances, the digraphs appear numerous times and because of that the size of the timing vector may differ from one digraph instance to another [27]. In a long text, there may be more than one instance of a digraph. However, the mean of all these instances is used as a corresponding latency time. It was shown in [28] that the typing pattern of a letter sequence may change when it is part of a larger word. For example, digraph ‘IS’ has different timing information in typing the word ‘**IS**’ and the word ‘**FUTURISTIC**’.

Figure 3: Digraphs and its number of occurrences.

Finally, we compute the confusion matrix in order to obtain the correct recognition rate. For each class, it presents the percentage of correctly classified users. We define our soft biometrics information as shown in Table 2.

Table 2: Soft Biometrics Information Class Label.

One Hand	=	1	Two Hands	=	-1
Male	=	1	Female	=	-1
< 30 years old	=	1	≥ 30 years old	=	-1
Right-handed	=	1	Left-handed	=	-1

4. Experimental Results

In the first subsection, we quantify the performance results of soft biometrics for keystroke dynamics with static passwords. As mentioned, distance measure are calculated for different timing information between two-

1
2
3
4
5
6
7
8
9 character sequences, and hence we show that with any combinations of two-
10 key characters (digraphs), significant results are obtained with free text il-
11 lustrated in the following subsection.
12
13
14

15 16 4.1. Passwords: Static

17
18 We performed several computations by using SVM. We recall that we
19 present the evolution of the average (over 100 computations) recognition
20 rate, while varying the percentage of the training data (from 1% to 90%), for
21 each of the four soft category.
22
23
24

25 26 • Hand Category Recognition

27
28 Figure 4 illustrates the results of the recognition rates for different training
29 ratios with one hand (C_1) and two hands (C_2) for five passphrases P_1 to P_5 .
30 To compute these results, an equal amount of data is used for both classes,
31 in particular 770 data samples for each class. In this experiment, the results
32 are promising, since from a ratio of training data over 50% of the total data
33 of the 110 users, the recognition rate is over 90%. This means that if the
34 database contains more than 55 users, the soft biometric system is able to
35 determine if the user types with one or two hands.
36
37
38
39
40
41
42
43

44 45 • Gender Category Recognition

46
47 Figure 5 illustrates the results of the recognition rates for different training
48 ratios with males (C_1) and females (C_2) for passphrases P_1 to P_5 . Only
49 30% of the data samples of male users are used (but all samples belong to
50 female) in order to have equilibrated classes (*i.e.* 224 data samples related
51 to male participants and 224 data samples related to female participants).
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9 The recognition rate depends on the particular passphrase and ranges from
10 70% to 86%.

11
12
13
14 • **Age Category Recognition**

15 Figure 6 illustrates the results of the recognition rates for different training
16 ratios with < 30 years old (C_1) and ≥ 30 years old (C_2) for passphrases P_1
17 to P_5 . We remove 46% of the data samples of class C_1 to have equal size
18 data classes each with the data of 51 users. The recognition rate for a ratio
19 over 50% is slightly less than the other soft criteria, namely between 67%
20 and 78%.

21
22
23
24
25
26
27
28 • **Handedness Category Recognition**

29 Figure 7 illustrates the results of the recognition rates for different training
30 ratios with right-handed (C_1) and left-handed (C_2) for passphrases P_1 to
31 P_5 . We keep only 12% of the right-handed class and all the left-handed
32 class to have equal size classes. The obtained recognition rate tends to vary
33 more than the other soft categories, but stays between 76% and 88%, which
34 are nevertheless quite good results. However, as mentioned, the selected
35 database for this category contains only 12 users in each class, therefore the
36 performances are decreased and the confidence intervals are wider compared
37 to other soft criteria with 110 users in each class.

38
39
40
41
42
43
44
45
46
47
48 • **Confidence Interval**

49 Table 3 shows the CI computed with a training dataset containing 50%
50 of the whole database, for different categories (*i.e.* hand, gender, age, hand-
51 edness).
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Figure 4: Average values for 100 iterations of recognition rates at 1% to 90% training ratios with two classes of hand for five passphrases [19].

Figure 5: Average values for 100 iterations of recognition rates at 1% to 90% training ratios with two classes of gender for five passphrases [19].

Figure 6: Average values for 100 iterations of recognition rates at 1% to 90% training ratios with two classes of age for five passphrases [19].

Figure 7: Average values for 100 iterations of recognition rates at 10% to 90% training ratios with two classes of handedness for five passphrases [19].

Table 3: Confidence interval computation at 50% training ratio for 5 passphrases and the data distribution (number of users) in each class.

	P_1	P_2	P_3	P_4	P_5
Hand (770 data samples: C_1&C_2)	96% \pm 0.1%	96% \pm 0.1%	95% \pm 0.1%	94% \pm 0.1%	94% \pm 0.1%
Gender (224 data samples: C_1&C_2)	74% \pm 0.3%	69% \pm 0.3%	70% \pm 0.2%	78% \pm 0.2%	76% \pm 0.2%
Age (357 data samples: C_1&C_2)	64% \pm 0.2%	64% \pm 0.2%	63% \pm 0.2%	69% \pm 0.2%	69% \pm 0.2%
Handedness (84 data samples: C_1&C_2)	72% \pm 1.2%	73% \pm 1.2%	72% \pm 1.2%	72% \pm 1.3%	73% \pm 1.2%

4.2. Free Text: Digraphs

We performed similar analysis with SVM as mentioned in Section 4.1. The first results deal with averaging of recognition rates (100 iterations) on all four soft categories for different percentage of training data ranging from 1% to 90%.

Figure 8 illustrates the evolution of the recognition rates on different training ratios with C_1 : one hand, male, age < 30 years old, right-handed; and C_2 : two hands, female, age \geq 30 years old for all four different soft biometrics information. In this experiment, the results are promising, from the ratio of 50% of total data used for SVM training, the recognition rate for Hand Class Recognition is over 90%; Gender Class Recognition is between 79% and 84%; Age Class Recognition is between 72% and 75%; and Handedness Class Recognition is between 83% and 88%. Table 4 summarises the performance comparison of recognition rate between passwords and free text at 50% to 90% training ratio.

Figure 8: Average values for 100 iterations of recognition rates at 1% to 90% training ratios with two classes of soft biometrics information with 14 digraphs (occurrences \geq 2) on free text.

4.3. Confusion Matrix: Majority Voting and Score Fusion for Passwords

In order to further enhance the performance, we perform data fusion, where we show that there is a great increase in the recognition accuracy rate results. The results of the obtained confusion matrix have improved significantly by fusing the data on all soft biometrics information at 50% training ratio based on *static passwords*. We apply the same equation and

Table 4: Summary of performance comparison of recognition rates for passwords and free text at 50% to 90% training ratios.

	Passwords	Free Text
Hand (770 data samples: C_1&C_2)	> 90%	> 90%
Gender (224 data samples: C_1&C_2)	between 70% and 86%	between 79% and 84%
Age (357 data samples: C_1&C_2)	between 67% and 78%	between 72% and 75%
Handedness (84 data samples: C_1&C_2)	between 78% and 88%	between 83% and 88%

ratio on free text, and here the results of the corresponding confusion matrix are based on *digraphs*. Then, the obtained performances are compared with three SVM computations: (i) before fusion for static passwords and free text; (ii) fusion based on *majority voting*; and (iii) fusion based on *score*; and only for static passwords in (ii) and (iii). Here, the fusion does not involve free text because all of the digraphs data are in the passwords. Table 5 summarises this information.

4.4. Discussions

From the previous results, we are able to see that the performances differ from one soft category to another. For static passwords, fusion processes namely *majority voting* and *score fusion* techniques have significantly increased the recognition performance rate on all soft biometrics characteristics from the initial results. *Score fusion*, however, gives better results, where we can see the results of this performance have increased significantly.

The results of free text are slightly superior to those of static passwords as illustrated in Table 5. As mentioned, we consider free text as the collection of the 5 passwords. With a total of only 14 occurrences consisting in three categories of digraph namely (i) 11 with two occurrences; (ii) 2 with three occurrences; and (iii) 1 with four occurrences, nevertheless, the results are quite promising.

Table 5: Performance comparison before and after fusion for passwords, and free text at 50% training ratio.

Technique	Soft Biometrics Information	Before Fusion	By fusing	
			Majority Voting	Score Fusion
Password	Hand Category	93.66%	100%	100%
	Gender Category	62.5%	85.71%	92.14%
	Age Category	55.49%	86.67%	85.71%
	Handedness Category	61.65%	84.52%	91.67%
Free Text	Hand Category	96.57%		
	Gender Category	80%		
	Age Category	65.71%		
	Handedness Category	83.33%		

5. Conclusions and Perspectives

In this paper, we propose a new soft biometric approach for keystroke dynamics. It consists of predicting the users' way of typing by defining the hand category *i.e.* number of hands used to type (one/two); gender category;

1
2
3
4
5
6
7
8
9 age category; and handedness category, where the results are promising for
10 both static passwords and free text. Moreover, we are able to enhance the soft
11 biometrics recognition rate for static passwords significantly by data fusion
12 and achieve higher performance accuracy. Another part of this work is a
13 comparative study between passwords and free text, where both approaches
14 give good results. For passwords, it is based on static texts where the users
15 are obliged to type with some constraints *i.e.* users type specific pre-defined
16 strings. Free text, on the other hand, with any combinations of two-key
17 characters (digraphs) is also able to provide good recognition rates without
18 any specific constraints on the users when typing.
19
20
21
22
23
24
25
26
27

28 For free text, it may be considered as suitable recognition if a set of
29 password is created by the user himself/herself, and hence having his/her
30 own freedom of texts choice. Nonetheless, the effectiveness of digraphs are
31 discriminative only when they are word-specific *i.e.* digraph features depend
32 on the word context they are occurred in [29]. Therefore, the obtained results
33 could be used as a reference model to assist the biometric system to better
34 recognise a user by a way he/she types on a keyboard. This will not only
35 strengthen the authentication process by hindering an impostor trying to
36 enter into the system, but also cut down on the computation time.
37
38
39
40
41
42
43
44

45 The results presented in this paper can be used to improve user authenti-
46 cation based on keystroke dynamics by combining two pieces of information:
47 (i) 'scores' provided by the biometric authentication system when comparing
48 the reference to a stored template; and (ii) a 'reliability index' by verifying
49 the concordance between one extracted soft biometric information (such as
50 gender) and the known information. The results in this work could also be
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9 applied, for example, in securing social networks, where the soft biometric
10 characteristics of a person in a chat can be checked against his/her claimed
11 profile.
12
13
14

15 16 **References**

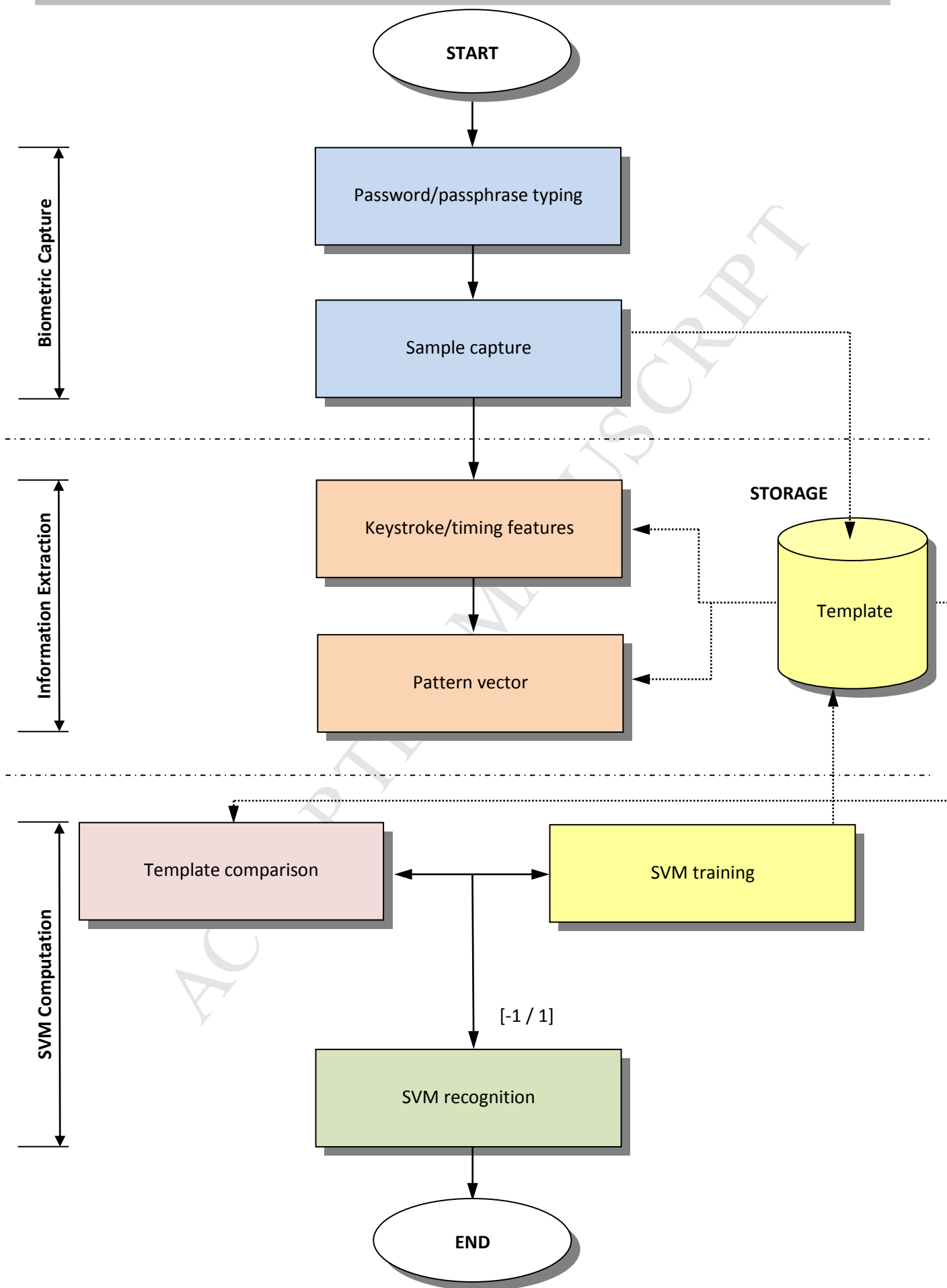
- 17
18
19 [1] P. Bours, Continuous keystroke dynamics: A different perspective
20 towards biometric evaluation, Information Security Technical Report
21 17 (1-2) (2012) p. 36–43. doi:10.1016/j.istr.2012.02.001.
22
23
24
25
26 [2] R. Giot, M. El-Abed, C. Rosenberger, Keystroke dynamics overview,
27 in: J. Yang (Ed.), Biometrics / Book 1, Vol. p. 1, InTech, 2011, pp. p.
28 157–182.
29
30 URL [http://www.intechopen.com/articles/show/title/
31 keystroke-dynamics-overview](http://www.intechopen.com/articles/show/title/keystroke-dynamics-overview)
32
33
34
35
36 [3] S. Impedovo, G. Pirlo, Verification of handwritten signatures: an
37 overview, in: 14th International Conference on Image Analysis and Pro-
38 cessing (ICIAP) 2007, IEEE, 2007, pp. p. 191–196.
39
40
41
42
43 [4] K. Moustakas, D. Tzovaras, G. Stavropoulos, Gait recognition using
44 geometric features and soft biometrics, Signal Processing Letters, IEEE
45 17 (4) (2010) p. 367–370.
46
47
48
49 [5] R. L. Klevans, R. D. Rodman, Voice recognition, Artech House, Inc.,
50 1997.
51
52
53
54 [6] F. Monrose, A. D. Rubin, Keystroke dynamics as a biometric for authen-
55 tication, Future Generation Computer Systems 16 (4) (2000) p. 351–359.
56
57
58

- 1
2
3
4
5
6
7
8
9 [7] D. Maio, A. K. Jain, Handbook of fingerprint recognition, Springer,
10 2009.
11
12
13 [8] R. Wildes, Iris recognition: an emerging biometric technology, Proceed-
14 ings of the IEEE 85 (9) (1997) p. 1348–1363.
15
16
17 [9] C. Epp, M. Lippold, R. Mandryk, Identifying emotional states using
18 keystroke dynamics, in: Proceedings of the 2011 Annual Conference on
19 Human Factors in Computing Systems, 2011, pp. p. 715–724.
20
21
22 [10] A. K. Jain, S. C. Dass, K. Nandakumar, Soft biometric traits for per-
23 sonal recognition systems, in: Proceedings of International Conference
24 on Biometric Authentication, Springer, 2004, pp. p. 731–738.
25
26
27 [11] H. Ailisto, E. Vildjiounaite, M. Lindholm, S.-M. Mkel, J. Peltola, Soft
28 biometrics—combining body weight and fat measurements with finger-
29 print biometrics, Pattern Recognition Letters 27 (5) (2006) p. 325 –
30 334.
31
32
33 [12] A. Dantcheva, C. Velardo, A. Dangelo, J.-L. Dugelay, Bag of soft biomet-
34 rics for person identification, Multimedia Tools and Applications 51 (2)
35 (2011) p. 739–777.
36
37
38 [13] H. T. F. Rhodes, Alphonse Bertillon, father of scientific detection,
39 Abelard-Schuman, 1956.
40
41
42 [14] U. Park, A. Jain, Face matching and retrieval using soft biometrics,
43 IEEE Transactions on Information Forensics and Security 5 (3) (2010)
44 p. 406 –415.
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4
5
6
7
8
9 [15] Y. Dong, D. L. Woodard, Eyebrow shape-based features for biometric
10 recognition and gender classification: A feasibility study, in: Interna-
11 tional Joint Conference on Biometrics (IJCB) 2011, IEEE, 2011, pp. p.
12 1–8.
13
14
15
16
17 [16] S. Denman, A. Bialkowski, C. Fookes, S. Sridharan, Determining op-
18 erational measures from multi-camera surveillance systems using soft
19 biometrics, in: 8th IEEE International Conference on Advanced Video
20 and Signal-Based Surveillance (AVSS) 2011, IEEE, 2011, pp. p. 462–467.
21
22
23
24
25 [17] G. L. Marcialis, F. Roli, D. Muntoni, Group-specific face verification
26 using soft biometrics, *Journal of Visual Languages & Computing* 20 (2)
27 (2009) p. 101–109.
28
29
30
31
32 [18] R. Bixler, S. D’Mello, Detecting boredom and engagement during writ-
33 ing with keystroke analysis, task appraisals, and stable traits, in: Pro-
34 ceedings of the 2013 international conference on Intelligent user inter-
35 faces, ACM, 2013, pp. p. 225–234.
36
37
38
39
40 [19] S. Idrus, E. Cherrier, C. Rosenberger, P. Bours, Soft biometrics for
41 keystroke dynamics, in: M. Kamel, A. Campilho (Eds.), *Image Analy-
42 sis and Recognition*, Vol. 7950 of *Lecture Notes in Computer Science*,
43 Springer Berlin Heidelberg, 2013, pp. p. 11–18.
44
45
46
47
48 [20] S. Idrus, E. Cherrier, C. Rosenberger, P. Bours, Soft biometrics
49 database: a benchmark for keystroke dynamics biometric systems, in:
50 2013 International Conference of the Biometrics Special Interest Group
51 (BIOSIG), 2013, pp. p. 281–288.
52
53
54
55
56
57
58

- 1
2
3
4
5
6
7
8
9 [21] V. Vapnik, Statistical learning theory, Wiley, 1998.
- 10
11 [22] C.-C. Chang, C.-J. Lin, Libsvm: a library for support vector machines,
12 ACM Transactions on Intelligent Systems and Technology (TIST) 2 (3)
13 (2011) p. 27, software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- 14
15
16
17
18
19
20 [23] C. Hsu, C. Chang, C. Lin, et al., A practical guide to support vector
21 classification (2003).
- 22
23
24 [24] M. Hearst, S. Dumais, E. Osman, J. Platt, B. Scholkopf, Support vector
25 machines, Intelligent Systems and their Applications, IEEE 13 (4) (1998)
26 p. 18–28.
- 27
28
29
30
31 [25] S. V. Stehman, Selecting and interpreting measures of thematic classifi-
32 cation accuracy, Remote sensing of Environment 62 (1) (1997) p. 77–89.
- 33
34
35 [26] R. Giot, M. El-Abed, B. Hemery, C. Rosenberger, Unconstrained
36 keystroke dynamics authentication with shared secret, Computers &
37 Security 30 (67) (2011) p. 427–445. doi:10.1016/j.cose.2011.03.004.
- 38
39
40
41
42 [27] H. Davoudi, E. Kabir, A new distance measure for free text keystroke
43 authentication, in: 14th International CSI Computer Conference (CS-
44 ICC) 2009, IEEE, 2009, pp. p. 570–575.
- 45
46
47
48
49 [28] T. Sim, R. Janakiraman, Are digraphs good for free-text keystroke dy-
50 namics?, in: Conference on Computer Vision and Pattern Recognition
51 (CVPR) 2007, IEEE, 2007, pp. p. 1–6.
- 52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4
5
6
7
8
9 [29] Y. Zhong, Y. Deng, A. K. Jain, Keystroke dynamics for user authentication, in: Conference on Computer Vision and Pattern Recognition Workshops (CVPRW) 2012, IEEE, 2012, pp. p. 117–123.
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

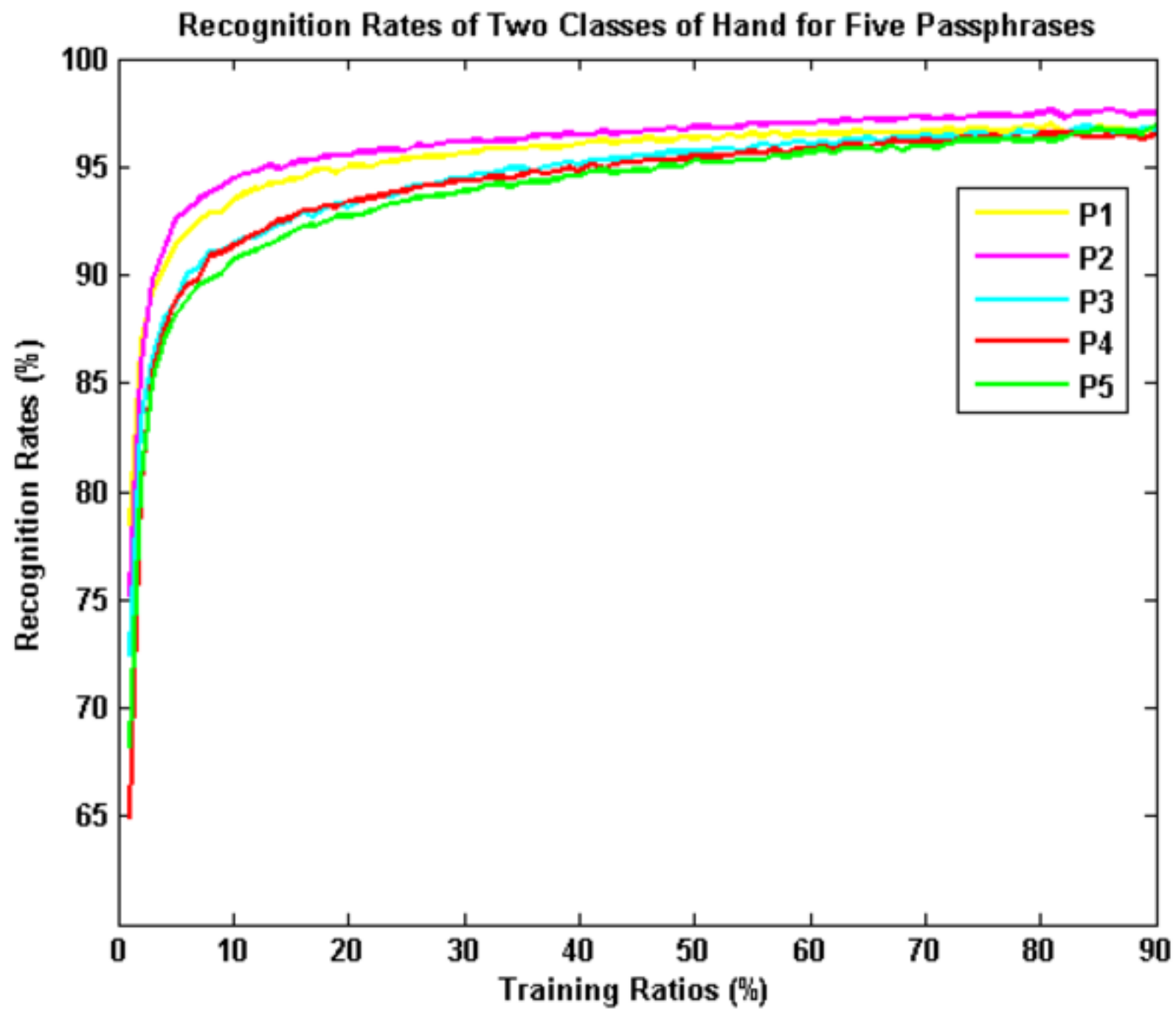


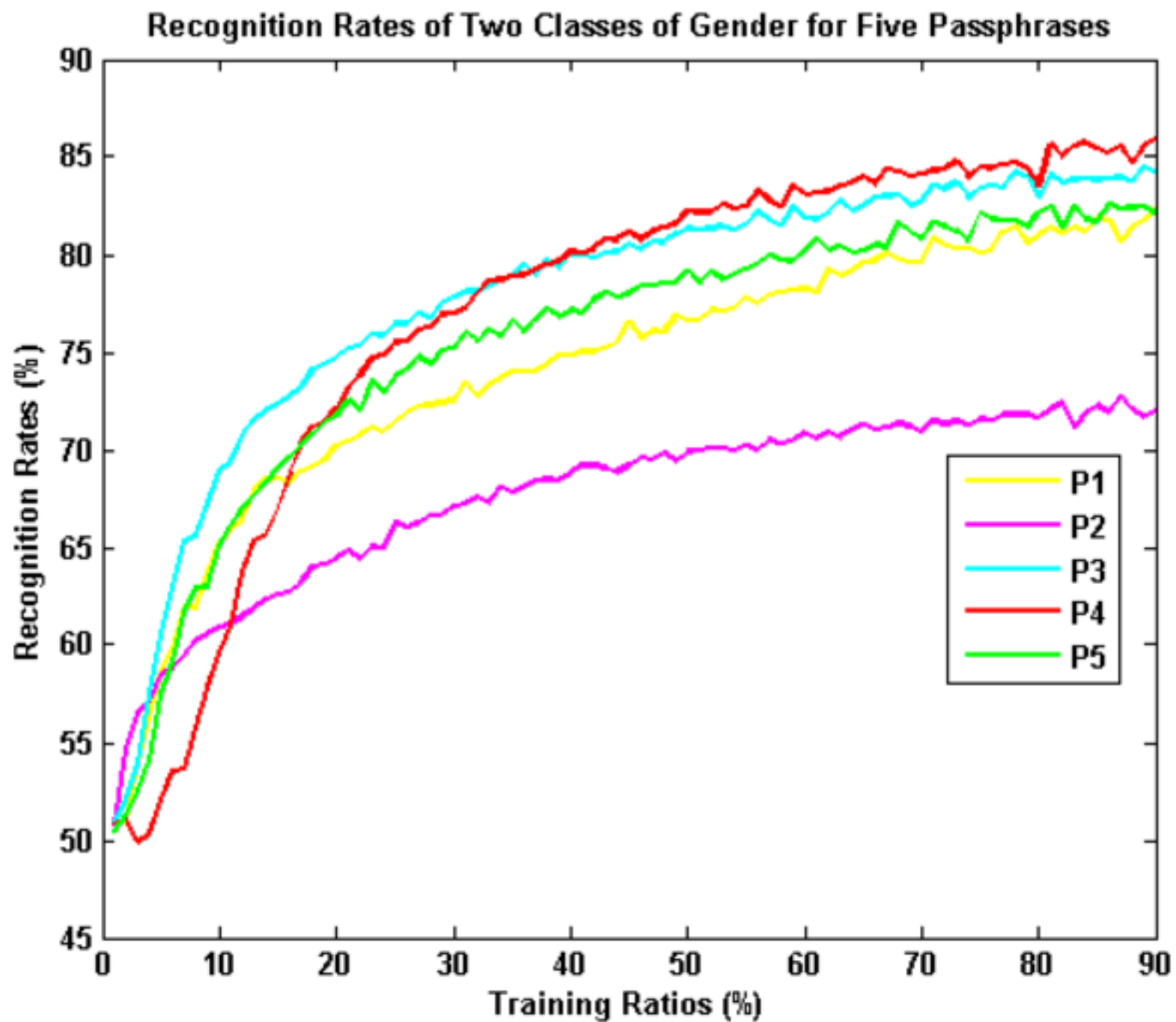
[Click here to download high resolution image](#)

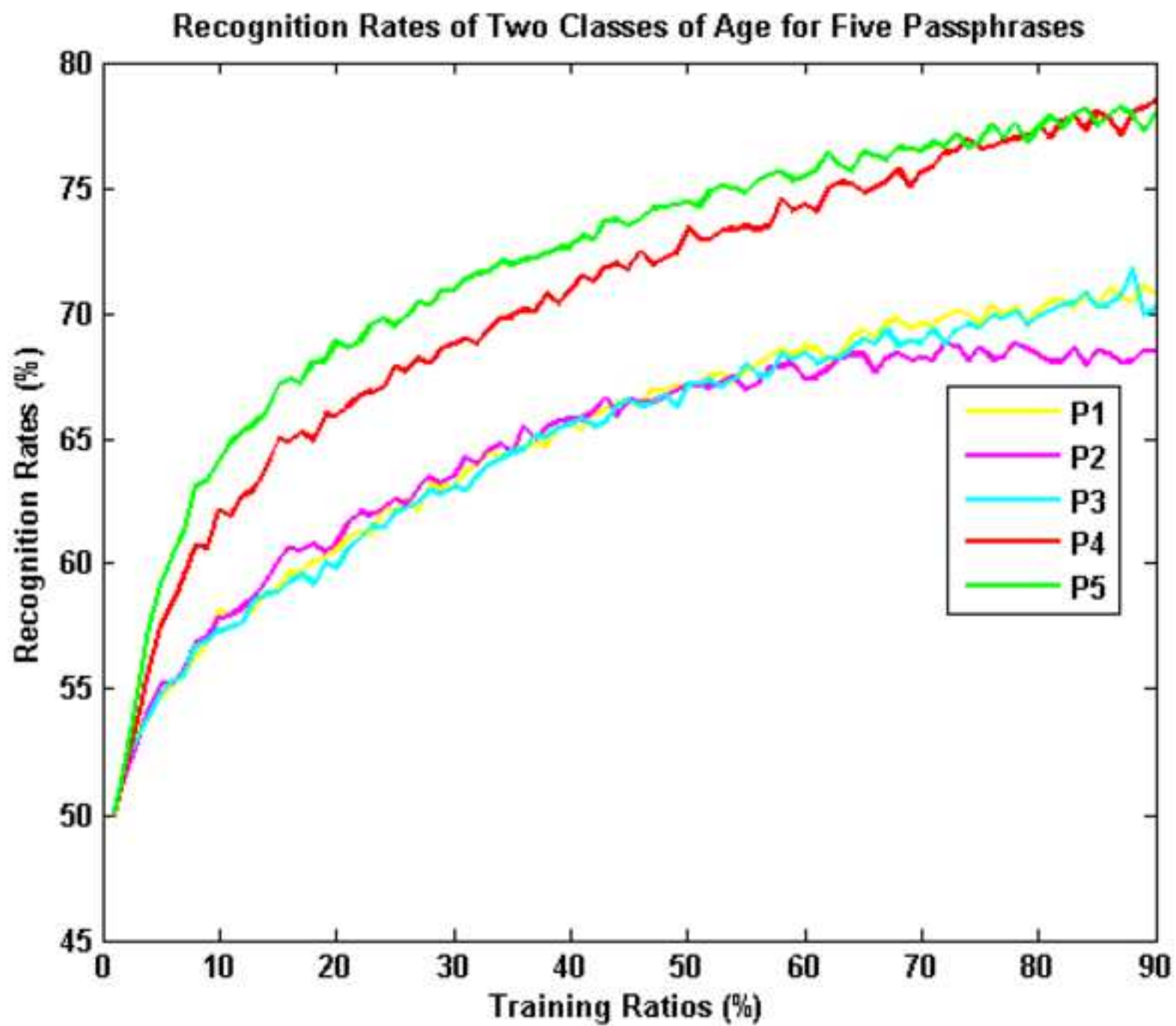
	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o	p	q	r	s	t	u	v	w	x	y	z	
a			■		■								■			■		■		■							
b																											
c	■■							■■■■																			
d									■						■						■						
e				■■								■			■			■■■	■■								
f																											
g																											
h	■				■■															■						■	
i			■■■									■		■	■						■						
j																											
k																											
l					■				■■		■■																
m	■				■				■																		
n	■				■	■			■																		
o						■								■■							■						
p					■■											■		■									
q																											
r				■	■				■■							■				■							
s			■																		■■						
t	■				■■			■								■											
u													■	■													
v																											
w																											
x																											
y																											
z																											

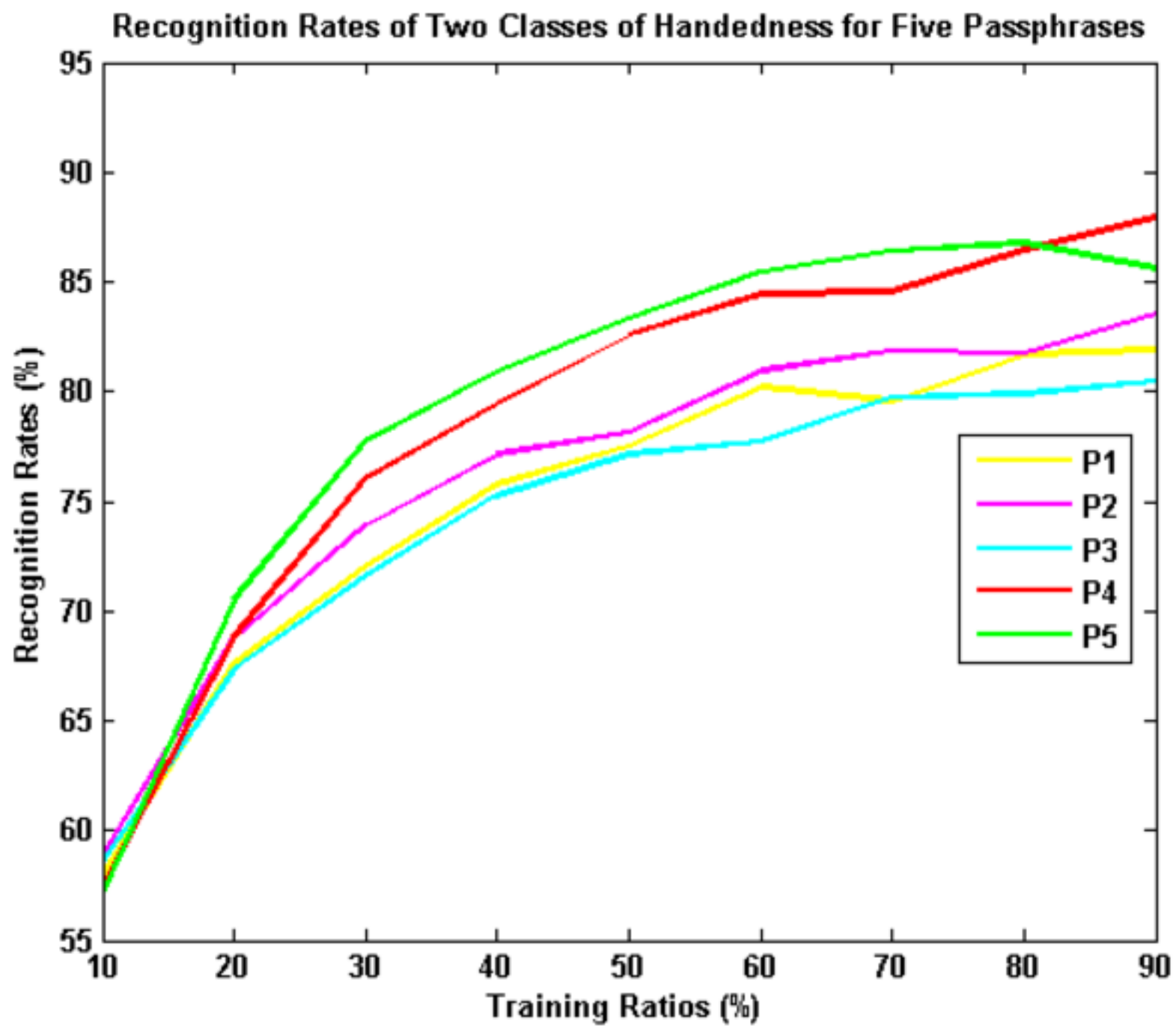
Passphrase:

- 1- leonardo dicaprio
- 2- the rolling stones
- 3- michael schumacher
- 4- red hot chilli peppers
- 5- united states of america

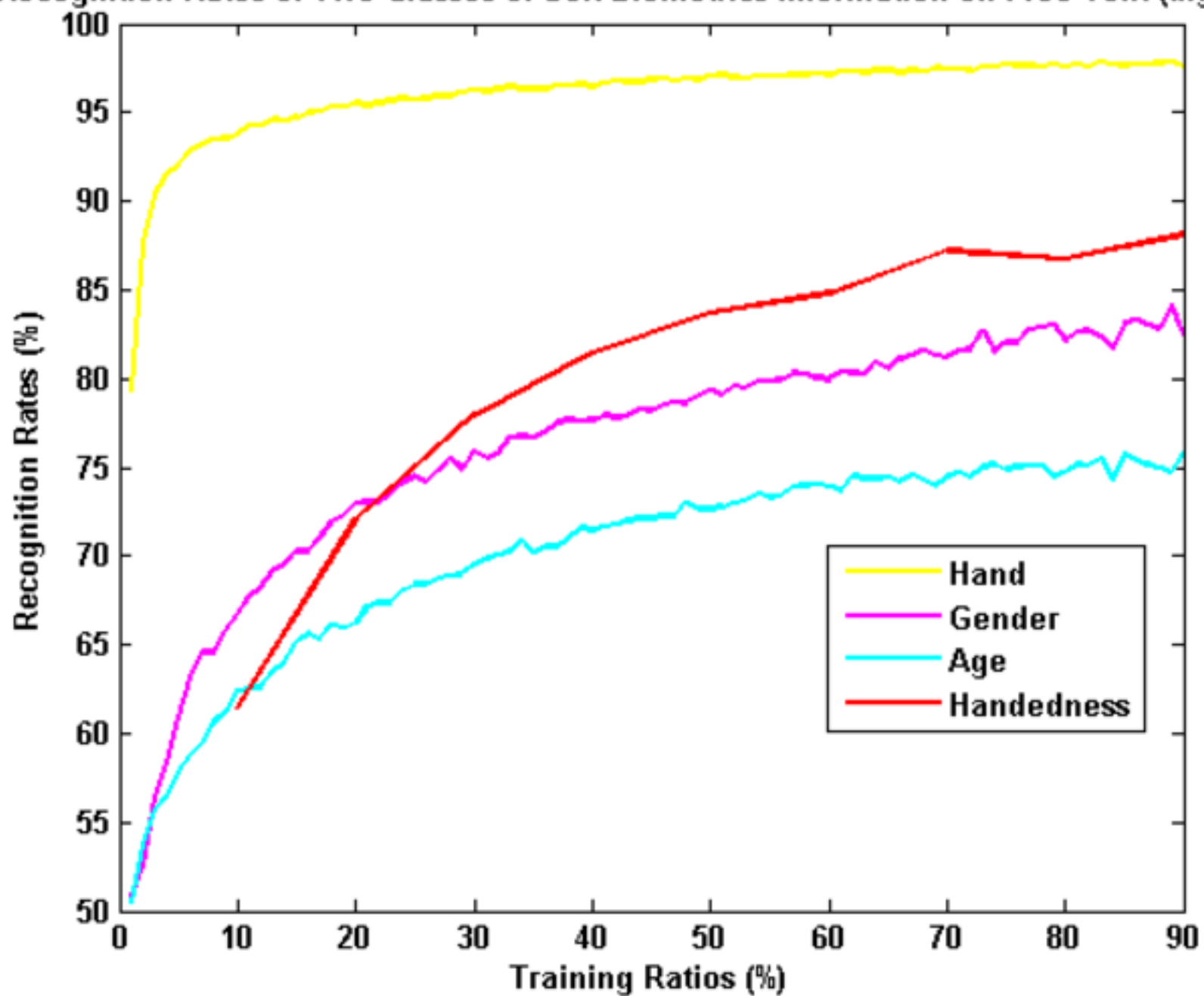








Recognition Rates of Two Classes of Soft Biometrics Information on Free Text (digraphs)



Authors Biographical Sketch for CoSe

Syed Zulkarnain Syed Idrus received the B.Sc. degree in Information Systems Engineering from University of Manchester Institute of Science and Technology (UMIST), United Kingdom and M.Sc. degree in Computer Engineering from Universiti Malaysia Perlis (UniMAP), Malaysia in 2001 and 2008, respectively. He is a Senior Lecturer at UniMAP (from 2009) and currently pursuing a Ph.D. degree in Computer Science and Applications at University of Caen Lower-Normandy, France specialising in biometrics. His research interest includes biometrics, pattern recognition, encryption, and information security.

Estelle Cherrier is an Associate Professor at ENSICAEN, France. She obtained her Ph.D. degree from the Collegium ingénieur de l'Université de Lorraine in 2006. She works at the GREYC Laboratory where she is a permanent member of the research group in E-payment & Biometrics. Her research interests include biometrics, signal processing and chaos system.

Christophe Rosenberger is a Full Professor at ENSICAEN, France. He obtained his Ph.D. degree from the University of Rennes I in 1999. He works at the GREYC Laboratory where he leads the research group in E-payment & Biometrics. His research interests include biometrics (definition of biometric systems and privacy issues). He is involved in developing authentication solutions for e-transactions applications.

Patrick Bours studied mathematics at the Eindhoven University of Technology in the Netherlands. He got his M.Sc. and Ph.D. with a specialisation in coding theory. After his studies, he worked at the Netherlands National Communication Security Agency (NLNCSA) in the area of cryptology. In July 2005, he moved to Norway where he joined the Norwegian Information Security Laboratory (NISlab), which is a part of Gjøvik University College. Since September 2009, he holds a Full Professor position. He is specialised in behavioural biometrics. His current research focus is on gait recognition, and static and continuous keystroke dynamics.