# ICDAR 2011 - Arabic Recognition Competition: Multi-font Multi-size Digitally Represented Text

Fouad Slimane[1,2] - Slim Kanoun[4] - Haikal El Abed[5] - Adel M. Alimi[2] - Rolf Ingold[1] - Jean Hennebert[1,3]

[1]*DIVA: Document, Image and Voice Analysis research group , Department of Informatics*
*University of Fribourg (unifr), Bd de Pérolles 90, CH-1700 Fribourg, Switzerland*
[2]*REGIM: REsearch Group on Intelligent Machines*
*University of Sfax, National School of Engineers (ENIS), BP 1173, Sfax, 3038, Tunisia*
[3]*Computer Science Department, EIAFR, HES-SO // Fribourg, Switzerland*
[4]*University of Sfax, National School of Engineers (ENIS), BP 1173, Sfax, 3038, Tunisia*
[5]*Institute for Communications Technology (IfN), Braunschweig, Germany*
*Fouad.Slimane@unifr.ch, slim.kanoun@ieee.org, elabed@tu-bs.de, Adel.Alimi@ieee.org,*
*Rolf.Ingold@unifr.ch, Jean.Hennebert@hefr.ch*

*Abstract*—This paper describes the Arabic Recognition Competition: Multi-font Multi-size Digitally Represented Text held in the context of the 11[th] International Conference on Document Analysis and Recognition (ICDAR2011), during September 18-21, 2011, Beijing, China. This first competition used the freely available Arabic Printed Text Image (APTI) database. Several research groups have started using the APTI database and this year, 2 groups with 3 systems are participating in the competition. The systems are compared using the recognition rates at the character and word levels. The systems were tested on one test dataset which is unknown to all participants (set 6 of APTI database). The systems are compared on the most important characteristic of classification systems, the recognition rate. A short description of the participating groups, their systems, the experimental setup, and the observed results are presented.

*Keywords*-APTI Database; Arabic; Recognition; Competition;

## I. INTRODUCTION AND MOTIVATION

Over the past ten years, Arabic recognition systems have achieved considerable improvements. The growing availability of benchmarking databases [1] [2] [3] and the organization of competitions [4] [5] [6], have contributed to systematic comparisons of various strategies for the benefit of their improvement. While handwritten Arabic tasks are rather well covered, few printed Arabic databases and competitions are available. So far, most of the printed Arabic systems have actually been benchmarked on private or small-scale databases, making their comparison rather difficult. To the best of our knowledge, the only free database for Arabic printed text is the Arabic Printed Text Image database (APTI) [1]. The most interesting characteristics of APTI are : very large set of images for significant benchmarking ($>$ 45 millions images), large lexicon of 113'248 words, multi-font, multi-size and single word images. Potentially less difficult than handwritten Arabic text recognition, APTI remains challenging due to

the variabilities induced by the different fonts and sizes that, in some cases, change drastically the distributions of observed features. APTI is typically related to OCR and "screen-based" OCR inputs where the user grab and crop a part of the computer screen.

This competition was organized by the DIVA (Document, Image and Voice Analysis) research group from the university of Fribourg, Switzerland in collaboration with the REGIM (REsearch Group on Intelligent Machines) from the National School of Engineers of Sfax, Tunisia and the group at the Institute of Communications Technolgy (IFN) of the Technical University of Braunshweig, Germany.
The competition was organized in a "blind" manner. The participants were asked to send an executable version of their recognizer to the organizers who, in turns, arrange to run the systems against an unseen set of data. The participants were able to train and tune their systems using the public parts of APTI.

The paper is organized as follows. Section 2 summarizes the main characteristics of APTI database. Section 3 is dedicated to the competition protocols. In section 4, we present the participating systems. Results are discussed in Section 5 and are followed by conclusions.

## II. THE APTI DATABASE

The APTI database was developed to promote the research and development of Arabic printed word recognition systems. Available from July 2009, APTI is freely distributed to the scientific community for benchmarking purposes [1]. At the time of writing this paper, 17 research groups have started using the APTI database.

The APTI database was created in low-resolution "72 dot/inch" with a lexicon of 113'284 different Arabic words

---

[1]http://diuf.unifr.ch/diva/APTI/
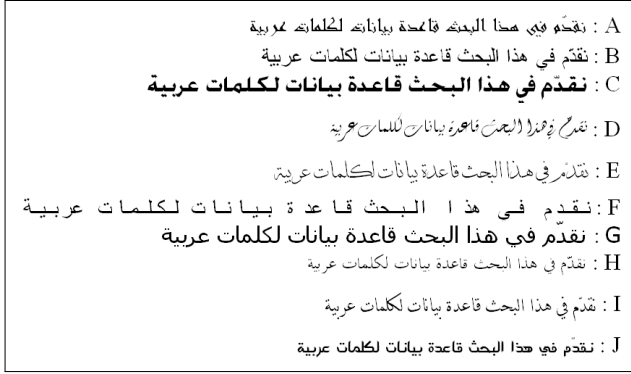
IEEE computer society

Figure 1. Fonts used to generate the APTI database: (A)Andalus, (B)Arabic Transparent, (C)AdvertisingBold, (D)Diwani Letter, (E)DecoType Thuluth, (F)Simplied Arabic, (G)Tahoma, (H)Traditional Arabic, (I)DecoType Naskh, (J)M Unicode Sara

```
<?xml version="1.0" encoding="UTF-8" ?>
- <wordImage id="78">
  - <content transcription="الآلاف" nPaws="4">
      <paw id="1" nbChars="1">Alif_I</paw>
      <paw id="2" nbChars="2">Laam_B TildAboveAlif_E</paw>
      <paw id="3" nbChars="2">Laam_B Alif_E</paw>
      <paw id="4" nbChars="1">Faa_I</paw>
    </content>
    <font name="Arabic Transparent" fontStyle="Plain" size="24" />
    <specs encoding="png" width="96" height="36" effect="none" />
    <generation type="downsampling5" renderer="java" filtering="antialiasing" />
  </wordImage>
```

Figure 2. Example of XML file including ground truth information about a given word image

of decomposable and non-decomposable words and 10 fonts presented in Figure 1. These fonts have been selected to cover different complexity of shapes of Arabic printed characters, going from simple fonts with no or few overlaps and ligatures (AdvertisingBold) to more complex fonts rich in overlaps, ligatures and flourishes (Diwani Letter). Different font sizes are also used in APTI: 6, 7, 8, 9, 10, 12, 14, 16, 18 and 24 points. We also used 4 different styles namely plain, italic, bold and combination of italic and bold. These sizes, fonts and styles are widely used on computer screen, Arabic newspapers and many other documents. The combination of fonts, styles and sizes guaranties a wide variability of images in the database. The total number of grey level images is above 45 million. Each word image in the APTI database is fully described using an XML file containing ground truth information about the sequence of characters as well as information about its generation (see an example on Figure 2). All Arabic letters have a natural distribution throughout the sets composing the database. The APTI lexicon includes 113'284 different single words. Table I shows the total quantity of word images, Piece of Arabic Words (PAWs) and characters in APTI.

The database is divided into six equilibrated sets to allow for flexibility in the composition of development and evaluation partitions. The words in each set are different

Table I
QUANTITY OF WORDS, PAWS AND CHARACTERS IN APTI

| | Words | PAWs | Characters |
|---|---|---|---|
| | 113'284 | 274'833 | 648'280 |
| | *10 Fonts * 10 Font Sizes * 4 Font Styles | | |
| Total | 45'313'600 | 109'933'200 | 259'312'000 |

but the distribution of all used letters is nearly the same in the various sets. For more details about statistics of each shape of characters in different sets, we refer to [1] and [7]. Images presented in table II shows some variabilities of APTI images thanks to the combination of fonts and sizes : they present the artefacts of the downsampling and antialiasing filters and the various forms of ligatures and overlaps of characters.

## III. THE COMPETITION

We invited groups working on Arabic word recognition to adapt their system to the APTI database and send us executables of their systems. The scientific objectives of this first edition are to measure the impact of font size on the recognition performances. This is evaluated in mono-font and multi-font contexts. The protocols are defined to evaluate the capacity of recognition systems to handle different sizes and fonts using digitally low resolution images in the aim to look for a robust approach to screen based OCR.

The evaluation has been organized using a blind procedure. The testing data of the evaluation is composed by an unpublished set (so called *set 6* of APTI) which is kept secret for evaluation purposes.

The evaluation will be reported as word and character recognition rates. In this first edition of the competition, we proposed 2 protocols:

1) First APTI Protocol for Competition: APTIPC1
   *Tested Fonts : Arabic Transparent*
   *Tested Style : Plain*
   *Tested Sizes : 6, 8, 10, 12, 18, 24*
   *Set 6 word images : 18'866 for each size/font*
   *Number of test in APTIPC1 : 6*

2) Second APTI Protocol for Competition: APTIPC2
   *Tested Fonts : Diwani letter, Andalus, Arabic Transparent, Simplified Arabic and Traditional Arabic*
   *Tested Style : Plain*
   *Tested Sizes : 6, 8, 10, 12, 18, 24*
   *Set 6 word images : 18'866 for each size/font*
   *Number of test in APTIPC2 : 30*

## IV. PARTICIPATING SYSTEMS

This section gives a short description of the submitted systems to the competition. The system descriptions vary in length due to the level of detail provided by the participants.

| Font/Size | 6 | 8 | 10 | 12 | 18 | 24 |
|---|---|---|---|---|---|---|
| *Andalus* | | | | | | |
| Image Resizing | 21X13 to 81X49 pixels | 28X17 to 81X49 pixels | 34X21 to 81X49 pixels | 41X25 to 81X49 pixels | 61X37 to 81X49 pixels | 81X49 to 81X49 pixels |
| *Arabic Transparent* | | | | | | |
| Image Resizing | 22X12 to 85X45 pixels | 29X16 to 85X45 pixels | 36X19 to 85X45 pixels | 43X23 to 85X45 pixels | 64X34 to 85X45 pixels | 85X45 to 85X45 pixels |
| *Simplified Arabic* | | | | | | |
| Image Resizing | 22X12 to 85X46 pixels | 29X16 to 85X46 pixels | 36X19 to 85X46 pixels | 43X23 to 85X46 pixels | 64X34 to 85X46 pixels | 85X46 to 85X46 pixels |
| *Traditional Arabic* | | | | | | |
| Image Resizing | 17X10 to 64X39 pixels | 22X13 to 64X39 pixels | 36X19 to 64X39 pixels | 43X23 to 64X39 pixels | 48X29 to 64X39 pixels | 64X39 to 64X39 pixels |
| *Diwani Letter* | | | | | | |
| Image Resizing | 21X12 to 80X44 pixels | 27X15 to 80X44 pixels | 34X19 to 80X44 pixels | 41X23 to 80X44 pixels | 60X33 to 80X44 pixels | 80X44 to 80X44 pixels |

## A. IPSAR System

IPSAR System was submitted by Samir Ouis, Mohammad S. Khorsheed and Khalid Alfaifi members in the Image Processing and Signal Analysis & Recognition (IPSAR) Group. This group is part of the Computer Research Institute (CRI) at King Abdulaziz City for Science & Technology (KACST) from the Saudi Arabia.

IPSARec is a cursive Arabic script recognition system where ligatures, overlaps and style variation pose challenges to the recognition system. It is based on Hidden Markov Model Toolkit (HTK). This is a portable toolkit for speech recognition system which is customized here to recognize characters. IPSARec is an omnifont, unlimited vocabulary recognition system. It does not require segmentation. The proposed system proceeds on three main stages: extracting a set of features from the input images, clustering the feature set according to a pre-defined codebook and finally, recognizing the characters.

Each word/line image is transferred into a sequence of feature vectors. Those features are extracted from overlapping vertical windows, divided into cells where each cell includes a predefined number of pixels, along the word/line image, then clustered into discrete symbols.

Stage two is performed within HTK. It couples the feature vectors with the corresponding ground truth to estimate the character model parameters. The final output of this stage is a lexicon-free system to recognize cursive Arabic text. During recognition, an input pattern of discrete symbols representing the word/line image is injected to the global model which outputs a stream of characters matching the text line.

For more details about this system, we refer to [8].

## B. UPV-BHMM Systems

These systems were submitted by Ihab Alkhoury, Adria Gimenez, and Alfons Juan, from the Universitat Politecnica de Valencia (UPV), Spain. They are based on Bernoulli HMMs (BHMMs), that is, HMMs in which conventional Gaussian mixture density functions are replaced with Bernoulli mixture probability functions [9]. Also, in contrast to the basic approach followed in [9], in which narrow, one-column slices of binary pixels are fed into BHMMs, the UPV-BHMM systems are based on a sliding window of adequate width to better capture image context at each horizontal position of the word image. This new, windowed version of the basic approach is described in [10]. As an example, Figure 3 shows the generation of a 7 X 5 word image of the number 31 from a sequence of 3 windowed (W = 3) BHMMs for the characters 3, "space" and 1.

The UPV-PRHLT systems were trained from input images scaled in height to 40 pixels (while keeping the aspect ratio) after adding a certain number of white pixel rows to both top and bottom sides of each image, and then binarized with the Otsu algorithm. A sliding window of width 9 was applied, and thus the resulting input (binary) feature vectors for the BHMMs had 360 bits. The number of states per character
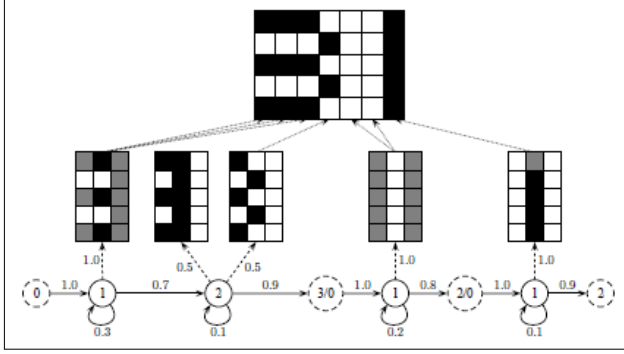
Figure 3. Generation of a 7 5 word image of the number 31 from a sequence of 3 windowed (W = 3) BHMMs for the characters 3, "space" and 1.

| System/Size | | 6 | 8 | 10 | 12 | 18 | 24 | Mean RR |
|---|---|---|---|---|---|---|---|---|
| IPSAR System | WRR | 5.7 | 73.3 | 75.0 | 83.1 | 77.1 | 77.5 | 65.3 |
| | CRR | 59.4 | 94.2 | 95.1 | 96.9 | 95.7 | **96.8** | 89.7 |
| UPV-PRHLT-REC1 | WRR | **94.5** | **97.4** | **96.7** | 92.5 | 84.6 | 84.4 | **91.7** |
| | CRR | **99.0** | **99.6** | **99.4** | 98.7 | 96.9 | 96.0 | **98.3** |
| UPV-PRHLT-REC2 | WRR | 94.5 | 97.4 | 96.7 | 92.5 | 84.6 | 84.4 | 91.7 |
| | CRR | 99.0 | 99.6 | 99.4 | 98.7 | 96.9 | 96.0 | 98.3 |
| DIVA-REGIM | WRR | 86.9 | 95.9 | 95.7 | 93.9 | 97.9 | 98.9 | 94.8 |
| | CRR | 98.0 | 99.2 | 99.3 | 98.8 | 99.7 | 99.7 | 99.1 |

was adjusted to 5 states for images with font size of 6, and 6 states for other font sizes. Similarly, the number of mixture components per state was empirically adjusted to 64. On the other hand, parameter estimation and recognition were carried out using the EM algorithm. Two systems were submitted: **UPV-PRHLT-REC1** and **UPV-PRHLT-REC2**. They are used for both tasks/protocols. In the first task (one style), there are no differences between systems, where one model for each font size is trained and used later to recognize the test corpus. For the second task: In the first system, for each font size, a different model for each font style is trained. The test corpus is recognized on all models, and recognized text word of the highest probability is selected. In the other system, a different character is considered for each style. A model for all styles together is trained and used to recognize the test corpus.

## V. TESTS AND RECOGNITION RESULTS

All systems have been tested using the *set 6* (18'866 single word images) of APTI database in different sizes and fonts. All participants sent us a running version of their recognition systems. The systems can be classified in two classes depending on the operating system: 2 systems are developed under Linux (UPV-PRHLT-REC1 and UPV-PRHLT-REC2) and one system under Microsoft Windows environment.

Table III presents all system results of the first APTI protocol (APTIPC1), we added the results of our system (DIVA-REGIM System) that we declared here "out of competition" for sake of integrity. The DIVA-REGIM system has actually been tuned for more than one year on set 1 to 5 of APTI. It is based on HMMs. One of its main characteristics is to be open vocabulary, i.e. able to recognize any Arabic printed word. The used HMM sub-models correspond to Arabic characters completed with a selected set of their corresponding variations, as explained in [11]. Similar character shapes are grouped into 65 models according to

the two following rules: (1) beginning and middle shapes share the same model, (2) isolated and end shapes share the same model. The feature vector is extracted from each analysis window. Using a simple right-left sliding procedure of the analysis window, no segmentation into letters is made and the word image is transformed into a sequence of feature vectors. During training time, the Expectation-Maximization (EM) algorithm is used to iteratively refines the component weights, means and variances to monotonically increase the likelihood of the training feature vectors. At recognition time, an ergodic HMM is built from all character models (i.e., all possible transitions between models are allowed).

For each test the best result is marked in bold. This first test is mono font and mono size. The test images presented to the systems are the one rendered using the font "Arabic Transparent", plain and sizes 6, 8, 10, 12, 18 and 24. For most of the systems, we observe good results in character recognition and slightly worse for the word recognition. both UPV-BHMM Systems have the same behaviour and show the best results with an average of 91.7 % for the word recognition rate and 98.3 % for the character recognition rate. Compared to other competition systems, the IPSAR system have the best character recognition rate on size 24.

Tables IV, V and VI presents system results of the second APTI protocol (APTIPC2) for competition. This second test is multi fonts and mono size. The test images presented to the systems are the one rendered using the fonts ("Arabic Transparent", "Andalus", "Simplified Arabic", "Traditional Arabic" and "Diwani Letter"), plain and sizes 6, 8, 10, 12, 18 and 24).

In the APTIPC2, the recognition rate is less good than in the APTIPC1 for the "Arabic Transparent" font. The best system is the UPV-PRHLT-REC1 with an average of 83.4 % for the word recognition rate and 96.4 % for the character recognition rate.

The UPV-PRHLT-REC1 system share good results for most fonts and sizes in this APTIPC2. The IPSAR system gives good results for the "Traditional Arabic" and "Diwani Letter" fonts in font size 10, 12 and 24.

Table IV
APTIPC2 - IPSAR SYSTEM RESULTS

| Font/Size | | 6 | 8 | 10 | 12 | 18 | 24 | Mean RR |
|---|---|---|---|---|---|---|---|---|
| Andalus | WRR | 13.9 | 35.7 | 65.6 | 73.8 | 69.5 | 64.5 | 53.8 |
| | CRR | 67.4 | 82.4 | 92.4 | 94.4 | 93.0 | 92.5 | 87.0 |
| Arabic Transparent | WRR | 29.9 | 40.0 | 73.2 | 74.9 | 65.9 | 69.1 | 58.8 |
| | CRR | 78.2 | 84.4 | 94.1 | 95.1 | 93.9 | 95.5 | 90.2 |
| Simplified Arabic | WRR | 30.8 | 39.8 | 73.2 | 75.5 | 66.2 | 68.6 | 59.0 |
| | CRR | 77.6 | 84.3 | 94.2 | 94.9 | 93.1 | 94.4 | 89.8 |
| Traditional Arabic | WRR | 4.6 | 3.4 | **46.7** | **55.1** | **52.9** | 50.4 | 35.5 |
| | CRR | 49.8 | 49.2 | **85.9** | **88.5** | **87.5** | 88.3 | 74.9 |
| Diwani Letter | WRR | 9.7 | 3.3 | **39.9** | **55.8** | 49.5 | **64.0** | 37.0 |
| | CRR | 60.1 | 48.3 | **83.4** | **89.1** | 91.7 | **92.6** | 77.5 |

Table V
APTIPC2 - UPV-PRHLT-REC1 SYSTEM RESULTS

| Font/Size | | 6 | 8 | 10 | 12 | 18 | 24 | Mean RR |
|---|---|---|---|---|---|---|---|---|
| Andalus | WRR | **94.1** | **75.5** | **81.1** | **83.6** | **83.9** | **85.0** | **83.8** |
| | CRR | **98.9** | **94.8** | **96.1** | **96.7** | **96.7** | **97.0** | **96.7** |
| Arabic Transparent | WRR | **94.7** | 78.2 | 78.9 | **81.8** | **83.1** | **83.8** | **83.4** |
| | CRR | **99.0** | 95.2 | 95.5 | 96.1 | **96.2** | **96.1** | **96.4** |
| Simplified Arabic | WRR | **95.8** | 82.4 | **84.2** | **85.3** | **85.6** | **88.0** | **86.9** |
| | CRR | **99.2** | 96.2 | **96.7** | **96.9** | **97.0** | **97.4** | **97.2** |
| Traditional Arabic | WRR | **57.6** | **38.3** | 43.6 | 43.5 | 42.9 | 46.2 | **45.4** |
| | CRR | **89.3** | **81.9** | 84.3 | 83.6 | 83.5 | 85.0 | **84.6** |
| Diwani Letter | WRR | **61.7** | **27.7** | 30.9 | 31.6 | **76.4** | 35.1 | **43.9** |
| | CRR | **90.9** | **75.8** | 77.8 | 78.1 | **94.9** | 79.6 | **82.8** |

## VI. CONCLUSIONS

APTI is challenging especially when we consider recognition rate at word level. We can observe that the impact of the character size is rather significant on the performances. Some systems perform better on larger size while other perform better on smaller size. This behaviour is probably conditioned by the training data and tuning for building the systems. The impact is even larger when systems have to deal with different fonts. For example, we can see that the word recognition performances are almost divides by 2 for most systems going from "Simplified Arabic" and "Diwani Letter".

The objective of the first competition for the recognition of multi-font and multi-size Arabic text was to evaluate and compare different systems and approaches. Two groups with 3 systems have participated at this first ICDAR 2011 Arabic Recognition Competition on digitally represented text. All systems are based on HMMs. The system UPV-PRHLT-REC1 is the winner of this first competition.

## VII. AKNOWLADGMENTS

## REFERENCES

[1] F. Slimane, R. Ingold, S. Kanoun, A. M. Alimi, and J. Hennebert, "A new arabic printed text image database and evaluation protocols," *Proc. of the 10th Int. Conf. on Doc. Analysis and Recognition (ICDAR)*, pp. 946–950, 2009.

[2] M. Pechwitz, S. Snoussi Maddouri, V. Margner, N. Ellouze, and H. Amiri, "Ifn/enit - database of handwritten arabic words," *CIFED*, pp. 129–136, 2002.

[3] Y. Al-Ohali, M. Cheriet, and C. Suen, "Databases for recognition of handwritten arabic cheques," *International Workshop on Frontiers in Handwriting Recognition*, pp. 601–606, 2000.

[4] H. E. Abed, V. Margner, M. Kherallah, and A. M. Alimi, "Icdar 2009 online arabic handwriting recognition competition," *Proc. of the 10th Int. Conf. on Doc. Analysis and Recognition (ICDAR)*, vol. 0, pp. 1388–1392, 2009.

[5] V. Margner and H. El Abed, "Arabic handwriting recognition competition," *Proc. of the 9th Int. Conf. on Doc. Analysis and Recognition (ICDAR)*, vol. 2, pp. 1274–1278, 2007.

[6] V. Margner and H. E. Abed, "Icdar 2009 arabic handwriting recognition competition," *Proc. of the 10th Int. Conf. on Doc. Analysis and Recognition (ICDAR)*, vol. 0, pp. 1383–1387, 2009.

[7] F. Slimane, R. Ingold, S. Kanoun, A. M. Alimi, and J. Hennebert, "Database and evaluation protocols for arabic printed text recognition," *DIUF-University of Fribourg - Switzerland*, 2009.

[8] M. S. Khorsheed, "Offline recognition of omnifont arabic text using the hmm toolkit (htk)," *Pattern Recogn. Lett.*, vol. 28, pp. 1563–1571, 2007.

[9] A. Gimenez and A. Juan, "Embedded bernoulli mixture hmms for handwritten word recognition," in *Proc. of the 10th Int. Conf. on Doc. Analysis and Recognition (ICDAR)*, 2009, pp. 896–900.

[10] A. Gime andnez, I. Khoury, and A. Juan, "Windowed bernoulli mixture hmms for arabic handwritten word recognition," in *2010 International Conference on Frontiers in Handwriting Recognition (ICFHR)*, 2010, pp. 533 –538.

[11] F. Slimane, R. Ingold, S. Kanoun, A. M. Alimi, and J. Hennebert, "Impact of character models choice on arabic text recognition performance," *International Conference on Frontiers in Handwriting Recognition(ICFHR)*, vol. 0, pp. 670–675, 2010.

Table VI
APTIPC2 - UPV-PRHLT-REC2 SYSTEM RESULTS

| Font/Size | | 6 | 8 | 10 | 12 | 18 | 24 | Mean RR |
|---|---|---|---|---|---|---|---|---|
| Andalus | WRR | 83.1 | 73.6 | 79.5 | 77.7 | 71.1 | 71.7 | 76.1 |
| | CRR | 96.0 | 94.1 | 95.1 | 94.9 | 93.6 | 93.5 | 94.5 |
| Arabic Transparent | WRR | 86.1 | **84.3** | **84.1** | 81.1 | 75.5 | 75.6 | 81.1 |
| | CRR | 97.1 | **96.5** | **96.6** | 96.1 | 94.9 | 94.8 | 96.0 |
| Simplified Arabic | WRR | 87.6 | **82.6** | 83.5 | 81.2 | 74.2 | 76.2 | 80.9 |
| | CRR | 97.4 | **96.1** | 96.5 | 96.1 | 94.7 | 95.0 | 96.0 |
| Traditional Arabic | WRR | 43.7 | 36.9 | 42.3 | 40.9 | 37.6 | 40.2 | 40.2 |
| | CRR | 83.6 | 80.5 | 83.2 | 82.1 | 80.8 | 82.2 | 82.1 |
| Diwani Letter | WRR | 41.9 | 26.4 | 29.7 | 29.2 | 68.4 | 29.9 | 37.6 |
| | CRR | 83.2 | 74.5 | 76.8 | 76.5 | 93.4 | 76.7 | 80.2 |