

UNIVERSITÉ DE FRIBOURG SUISSE UNIVERSITÄT FREIBURG SCHWEIZ

Database and Evaluation Protocols for Arabic Printed Text Recognition

Fouad Slimane^{1,2} Rolf Ingold¹ Slim Kanoun³ Adel M. Alimi² Jean Hennebert^{1,4}

{Fouad.Slimane, Jean.Hennebert, Rolf.Ingold}@unifr.ch, { Slim.Kanoun, Adel.Alimi}@enis.rnu.tn

> February 15, 2012 Version 1.1

DEPARTMENT OF INFORMATICS RESEARCH REPORT

Département d'Informatique – Department für Informatik • Université de Fribourg – Universität Freiburg • Boulevard de Pérolles 90 • 1700 Fribourg • Switzerland

Phone +41 (26) 300 84 65 fax +41 (26) 300 97 26 Diuf-secr-pe@unifr.ch http://diuf.unifr.ch

¹ DIVA-DIUF, University of Fribourg, Switzerland

² REGIM, University of Sfax, Tunisia

³ National School of Engineers, University of Sfax, Tunisia

⁴ HES-SO // Wallis, University of Applied Sciences Western Switzerland, Switzerland

Abstract

We report on the creation of a database composed of images of Arabic Printed Text. The purpose of this database is the large-scale benchmarking of open-vocabulary, multi-font, multi-size and multi-style text recognition systems in Arabic. Such systems take as input a text image and compute as output a character string corresponding to the text included in the image. The database is called APTI for Arabic Printed Text Image. The challenges that are addressed by the database are in the variability of the sizes, fonts and style used to generate the images. A focus is also given on low-resolution images where anti-aliasing is generating noise on the characters to recognize. The database is synthetically generated using a lexicon of 113'284 words, 10 Arabic fonts, 10 font sizes and 4 font styles. The database contains 45'313'600 single word images totaling to more than 250 million characters. Ground truth annotation is provided for each image thanks to a XML file. The annotation includes the number of characters, the number of PAWs (Pieces of Arabic Word), the sequence of characters, the size, the font used to generate each image, etc. The database is called APTI for Arabic Printed Text Images.

Keywords: Arabic Text Recognition system, benchmarking, text image databases, OCR

1 Introduction and motivations

With a quite large user base of about 300 million people worldwide, Arabic is important in the culture of many people. In the last fifteen years, most of the efforts in Arabic text recognition have been put for the recognition of scanned off-line printed documents [Khorsheed 07] [Husni 08] [Shaaban 08] [Slimane 08]. Most of these developments have been benchmarked on private databases and therefore, the comparison of systems is rather difficult.

To our knowledge, there are currently few large-scale image databases of Arabic printed text available for the scientific community. One of the only references we have found is about the ERIM database containing 750 scanned pages collected from Arabic books and magazines [Schlosser 95]. However, it seems difficult to have access to this database. In the field of Arabic handwriting recognition, public databases do exist such as the freely available IFN/ENIT-database [Pechwitz 02] Open competitions are even regularly organized using this database [Margner 05] [Margner 07].

On the other hand, text corpus or lexical databases in Arabic are available from different associations or institutes [Graff 06] [Abbes 04] [AbdelRaouf 08]. However, such text corpora are not directly usable for the benchmarking of recognition systems that take images as input.

Considering this, we have initiated the development of a large database of images of printed Arabic words. This database will be used for our own research and will be made available for the scientific community to evaluate their recognition systems. The database has been named APTI for Arabic Printed Text Image.

The purpose of the APTI database is the large-scale benchmarking of open-vocabulary, multi-font, multi-size and multi-style text recognition systems in Arabic. The images in the database are synthetically generated from a large corpus using automated procedures. The challenges that are addressed by the database are in the variability of the sizes, fonts and style used to generate the images. A focus is also given on low-resolution images where antialiasing is generating noise on the characters to recognize. By nature, APTI is well suited for the evaluation of screen-based OCR systems that take as input images extracted from screen captures or pdf documents. Performances of classical scanned-based OCR or camera-based OCR systems could also be measured using APTI. However, such evaluations should take into account the absence of typical artefacts present in scanned or camera documents.

While synthetically generated, the challenges of the database remain multiple:

- Large-scale evaluation with a realistic sampling of most of the Arabic character shapes and their accompanying variations due to ligatures and overlaps;
- Availability of multiple fonts, styles and sizes that must be nowadays treated by recognition systems;
- Emphasis on low resolution images that are nowadays frequently present on computer screens;
- Isolated word images where inter-word language models cannot be used;
- Semi-blind evaluation protocols with decoupled development/evaluation sets.

The objective of this paper is to describe the APTI database and the evaluation protocols defined on the database. In section 2, we present details about lexicon, fonts, font-sizes, rendering procedure, Sources of variability and ground truth description. In section 3, statistical information about the content of the database are given. The evaluation protocols are showed in section 4. Finally, some conclusions are presented in section 5.

2 Specifications of APTI-Database

The APTI database is synthetic and images are generated using automated procedures. In this section, we present the specification of this database.

2.1 Lexicon

The APTI database contains a mix of decomposable and non-decomposable words images. Decomposable words are generated from root Arabic verbs using Arabic schemes [Kanoun 2005] and non-decomposable words are formed by Arabic proper names, general names, country/town/village names, Arabic prepositions, etc.

To generate the lexicon, we have parsed different Arabic books such as *The Muqaddimah* - *An introduction to history of Ibn Khaldun⁵ and Al-bukhala of Gahiz⁶* as well as Arabic articles taken from the Internet. This parsing procedure totalled 113'284 single different Arabic words, leading to a pretty good coverage of the Arabic words mostly used in texts. The language used in our sources is exclusively in standard Arabic with no dialect.

2.2 Fonts, styles and sizes

Taking as input the words in the lexicon, the images of APTI are generated using 10 different fonts presented in Fig. 1: Andalus, Arabic Transparent, AdvertisingBold, Diwani Letter, DecoType Thuluth, Simplified Arabic, Tahoma, Traditional Arabic, DecoType Naskh, M Unicode Sara. 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 or Thuluth).

⁵ Ibn Khaldoun, (May 27,1332 – March 19, 1406) was a famous historien, scholar, theologian, and statesman born in North Africa in presentday Tunisia. (http://en.wikipedia.org/wiki/Ibn_Khaldoun)

⁶ Al-Jahiz, (born in Basra, c. 781 – December 868 or January 869) was a famous Arab scholar, believed to have been an Afro-Arab of East African descent.(http://en.wikipedia.org/wiki/Al-Jahiz)

Different sizes are also used in APTI: 6 points, 7 points, 8 points, 9 points, 10 points, 12 points, 14 points, 16 points, 18 points 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, books and many other documents. The combination of fonts, styles and sizes guaranties a wide variability of images in the database.

Overall, the APTI database contains 45'313'600 single words images, taking into account the full lexicon where the different combinations of fonts, style and sizes are applied.

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

2.4 Rendering procedure

The text images are generated using automated procedures. As a consequence, artefacts or noise usually present for scanned or camera-based documents are not present in the images. Such degradations could actually be artificially added, if needed [Baird 08], but it is currently out of the scope of APTI.

Image generation of text, for example on screen, can be done in many different ways. They are usually all leading to slight variations of the target image. We have opted for a rendering procedure that allows us to include effects of downsampling and antialiasing. These effects are interesting in terms of variability of the images, especially in low-resolution.

The procedure involves the downsampling of a high resolution source image into a low resolution image using antialiasing filtering. We also use different grid alignments to introduce variability in the application of the antialiasing filter. The details of the procedure are the following:

- 1. A gray-scale source image is generated in high resolution (360 pixels/inch) from the current word in the lexicon, using the selected font, size and style (Example in Fig. 2, height of image = 119, width of image =247).
- 2. Columns and rows of white pixels are added to the right hand side and to the top of the image. The number of columns and rows is chosen to have a height and width multiple of the downsampling factor (for example image in Fig. 3, we add 3 white columns and 1 white row). This effect allows to have the same deformation in all images and artificially moving the downsampling grid.
- 3. Downsampling and antialiasing filtering are applied to obtain the target image in lower resolution (72 pixels/inch) (Example in Fig. 3, height of image =24, width of image = 50). The target image is in grey level. The downsampling and antialiasing algorithms are the one implemented in the Java abstract class Image. In our implementation, we used the SCALE_SMOOTH option of the Java method which optimize the downsampling algorithm selection according to quality and speed.



Fig. 2: Example of Arabic image word source



Fig. 3: Example of anti-aliasing effect and down sampling result approach

In Fig. 3, Character "Alif" is presented in two different forms (different presentation of anti-aliasing effect) in the same word image although it has the same characteristics (Font, Font Size, Style,...).

2.5 Sources of variability

The sources of variability in the generation procedure of text images in APTI are the following:

- 1. 10 different fonts: Andalus, Arabic Transparent, AdvertisingBold, Diwani Letter, DecoType Thuluth, Simplified Arabic, Tahoma, Traditional Arabic, DecoType Naskh, M Unicode Sara;
- 2. 10 different sizes: 6, 7, 8, 9, 10, 12, 14, 16, 18 and 24 points;
- 3. 4 different styles: plain, bold, italic, italic and bold;
- 4. Various forms of ligatures and overlaps of characters thanks to the large combination of characters in the lexicon and thanks to the used fonts;
- 5. Very large vocabulary that allows to test systems on unseen data;
- 6. Various artefacts of the downsampling and antialiasing filters due to the random insertion of columns of white pixels at the beginning of image words;
- 7. Variability of the height of each word image.

The last point of the previous list is actually intrinsic to the sequence of characters appearing in the word. In APTI, there is actually no a priori knowledge of the position of the baseline and it is up to the recognition algorithm to compute the baseline, if needed.

2.6 Ground truth description

Each word image in APTI database is fully described using an XML file containing ground truth information about the sequence of characters as well as information about the generation. An example of such XML file is given in Fig. 4.

Fig. 4: Example of XML file including ground truth information about a given word image

The XML file is composed by four principal markups sections:

- *Content*: in this element, we have the transcription of Arabic word, the number of Piece of Arabic Word (nPaws) and sub-elements for each PAW with the sequence of characters. In our representation, characters are identified using plain English labels as described below.
- *Font*: in this element, we specify the font name, font style and size used to generate the image word.

- *Specs*: in this element, we present the encoding of image, width, height and eventual addition effect. In the current version of APTI, there is actually no added effects but we have planned to use this attribute for later versions of image rendering where effects could be present.
- *Generation*: in this element, we indicate the type of generation, the tool used for generation and the used filter in generation. In the current version of APTI, this element is constant as the same generation procedure has been applied. The type 'downsampling5' is here indicating that the generation procedure correspond to a downsampling, using factor 5, from high resolution source images as explained in Section 2.4.

Letter label	Number of Occurence	Isolate Begin		Middle	End
Alif	90353	1	-	L	-
Baa	28119	ب	÷	÷	÷
Тааа	59343	ت	ï	ž	Ľ
Thaa	3803	ث	ć	ن	ث
Jiim	11455	5	÷	ج	ج
Нааа	17866	ζ	~	ح	で
Xaa	8492	ċ	ż	خ	خ
Daal	18399	د		ح	
Thaal	3100	ذ		ذ	
Raa	37571	ر		بر	
Zaay	6325	ز		ز	
Siin	21648	س	ىد	ىىد	س
Shiin	8668	ش	ىثد	ىىتىر	ىش
Saad	8310	ص	صد	<u>مد</u>	ص
Daad	5548	ض	ضد	ضر	ض
Thaaa	8610	ط	ط	Ь	ط
Таа	1438	ظ	ظ	ᄨ	ظ
Ayn	16552	٤	ع	£	ع
Ghayn	5912	È	à	غ	غ
Faa	13749	ف	ف	غ	ف
Gaaf	16819	ق	ē	ā	ق
Kaaf	12711	ك	ک	ک	1ك
Laam	41159	J	J	T	ل
Miim	47084	م	\$	4	R
Nuun	44186	ن	ذ	ì	ىن
NuunChadda	1343	ڹۨ	ě	ž	ڹۨ
Наа	16094	٥	ھ	÷	4
Waaw	26008	9		و	
Yaa	40215	ى	2	*	J
YaaChadda	4348	يّ	2		ي
Hamza	1142	¢			
HamzaAboveAlif	8770	ĺ		Ĺ	
TaaaClosed	8376	ة			Ä
HamzaUnderAlif	1501	ļ		Ļ	

AlifBroken	972	ى			J
TildAboveAlif	500		ĩ		ĩ
HamzaAboveAlifBroken	1253	ئ	ć	د •	ئ
HamzaAboveWaaw	538	ۇ		ڂ	
Quantity of Characters	648'280				
Quantity of PAWs	274833				
Quantity of words	113'284				

Table 1: Arabic letters with used labels and occurrence in APTI database

The different character labels are summarized in Table 1. As the shape of characters are varying according to their position in the word, the character labels also include a suffix to specify the position of the character in the word: "B" standing for beginning, "M" for Middle, "E" for end and "I" for isolated. The character "Hamza" being always isolated, we don't use the position suffix for this character. We also artificially inserted characters labels such as "NuunChadda" or "YaaChadda" to represent the character shape issued from the combination of "Nuun" and "Chadda" or "Yaa" and "Chadda".

3 Database statistics

The APTI database consists of 113'284 different single words presented in 10 fonts, 10 font-sizes and 4 font-styles. Table 2 shows the total quantity of word images, PAWs (Piece of Arabic Words), and characters in APTI database.

	Number of Words	Number of PAWs	Number of characters	
	113'284	274'833	648'280	
Number of Font	10	10	10	
Number of Font	10	10	10	
Size	10	10	10	
Number of Font	1	1	1	
Styles	4	4	4	
Total	45'313'600	109'933'200	259'312'000	

Table 2: Quantity of words, PAWs, characters in database

3.1 Division into sets

We have divided the database into six equilibrated sets to allow for flexibility in the composition of development and evaluation partitions. The words in each set are different but the distribution of all used letters is nearly the same in the various sets (see Table 3). The five first sets are available for the scientific community and the sixth set is kept internal for potential future evaluation of systems in blind mode.

The algorithm for the distribution of words in the different sets has been designed to have similar allocations of letters and words in all sets. The algorithm is presented in details in Fig. 5. The steps of the algorithms are the following. First (step 1 in Fig. 5), we read all the words from the database and we accumulate the number of occurrence of each used letters. The letters are then sorted according to their number of occurrence, from the smallest number of occurrence to the largest. Second, (step 2 in Fig. 5), bins (vectors) are created for each letters and they are ordered according to the occurrences computed in step 1. For each word of the database, we go through the bins and we look if the word contains the character associated to the bin. If yes, then the word is associated to the bin and we go to the next word. Doing this,

we actually build sets of words having letters with low occurrences. Third (step 3 in Fig. 5), we go through each bin and distribute the word sequentially in our final 6 sets, emptying each bin one after the other.

In short, this procedure is simply stressing a fair distribution of words that include characters with few occurrences. Such a distribution is important to avoid that a given character is under-represented in a given set and therefore to avoid potential problems of during training or testing time.

```
# Inputs: list of Arabic Words
# Output: six Sets of Arabic words with similar distribution of words and characters
Begin
1. for all words w_i, i \in \{1...113'285\} in APTI
for all used letters l_i, i \in \{1...38\}
integer tab[]=findNbOccureneOfLetters(I<sub>i</sub>);
endfor
endfor
increasingSort(tab);
2. for all words w_i, i \in \{1...113'285\} in APTI
   for all used letters l_i, j \in \{1...38\} sorted by NbOccurence
         if (1_i \subset w_i)
                   add w<sub>i</sub> in vector V_i, j \in \{1...38\}
                   go to 2
         endif
    endfor
    endfor
3. for all V_s, s \in \{1...38\}
         read w_i, i \in \{1...NbWordInV_s\} from V_s
         if i mod 6=0
                   add w<sub>i</sub> in S<sub>1</sub>
         if i mod 6=1
                   add w_i in S_2
         if i mod 6=2
                   add w<sub>i</sub> in S<sub>3</sub>
         if i mod 6=3
                   add w_i in S_4
         if i mod 6=4
                   add w<sub>i</sub> in S<sub>5</sub>
         if i mod 6=5
                   add w<sub>i</sub> in S<sub>6</sub>
    endfor
end
```

Fig. 5: Algorithm used for distribution

Letter label	Set 1	Set 2	Set 3	Set 4	Set 5	Set 6
Alif	15078	14925	15165	15120	15046	15019
Baa	4513	4763	4692	4704	4730	4717
Тааа	9926	9884	9897	9797	9942	9897
Thaa	634	633	631	634	643	628
Jiim	1893	1897	1887	1924	1915	1939
Нааа	2953	2963	3017	2933	3000	3000
Xaa	1407	1435	1439	1401	1403	1407
Daal	3187	3033	3075	2990	3028	3086
Thaal	514	520	528	504	516	518
Raa	6304	6243	6169	6335	6253	6267
Zaay	1064	1054	1054	1066	1042	1045
Siin	3674	3556	3674	3512	3629	3603
Shiin	1457	1446	1418	1434	1455	1458
Saad	1374	1377	1388	1411	1371	1389
Daad	922	943	936	906	921	920
Thaaa	1419	1426	1431	1426	1446	1462
Таа	242	238	240	238	239	241
Ayn	2764	2823	2769	2718	2755	2723
Ghayn	981	970	983	984	990	1004
Faa	2305	2256	2221	2313	2339	2315
Gaaf	2784	2734	2853	2883	2762	2803
Kaaf	2101	2090	2099	2145	2136	2140
Laam	6745	6926	6972	7002	6790	6724
Miim	7871	7836	7957	7806	7797	7817
Nuun	7484	7433	7289	7316	7400	7264
NuunChadda	225	224	224	223	224	223
Наа	2670	2687	2590	2718	2705	2724
Waaw	4421	4313	4325	4333	4264	4352
Yaa	6641	6630	6876	6685	6648	6735
YaaChadda	725	727	709	719	735	733
Hamza	192	187	190	193	192	188
HamzaAboveAlif	1437	1483	1455	1512	1456	1427
TaaaClosed	1417	1407	1394	1364	1409	1385
HamzaUnderAlif	253	250	256	247	248	247
AlifBroken	162	161	164	163	161	161
TildAboveAlif	84	84	83	83	83	83
HamzaAboveAlifBroken	210	208	208	209	208	210
HamzaAboveWaaw	89	90	89	91	89	90
Quantity of Characters	108'122	107'855	<i>108'347</i>	108'042	107'970	107'944
Quantity of PAWs	45'982	45'740	45'792	45'884	45'630	45'805
Quantity of words	18897	18892	18886	18875	18868	18866

Table 3: Distribution of characters in the different sets of APTI

3.2 Distribution of letters in sets

Tables 4 to 9 are presenting the distribution of each shape of characters in their respective sets.

Letter label	Nb. Occ	Isolate	Begin	Middle	End
Alif	15078	15	1 5823		.55
Baa	4513	128 ب	1 978 .	4 2226	181 🕰
Тааа	9926	587 ت	3626	: 5332	381 ت
Thaa	634	12 ث	\$ 261	341	20 ث
Jiim	1893	ट 60	→ 781	× 1016	7. 36
Haaa	2953	7 69	> 1135	× 1648	7. 101
Xaa	1407	τ 16	.→ 587	→ 782	<u>22</u> ح
Daal	3187	و د	988	21 د	199
Thaal	514	: ذ	167	3 ند	47
Raa	6304	1 ر	813	44 ر	91
Zaay	1064	;	389	; 6	75
Siin	3674	68 س	ىىد	ىىد	89 س
			1434	2083	
Shiin	1457	18 ش	580 شد	831 مثد	28 ـش
Saad	1374	14 ص	439 صد	882 صد	39 ص
Daad	922	41 ض	358 ضد	497 ض	26 ض
Thaaa	1419	42 ط	392 ط	920 ط	65 ط
Таа	242	6 ظ	58 ظ	163 ظ	15 ط
Ayn	2764	E 67	1003 ع	s 1575	119ع
Ghayn	981	<u>1</u> 2 غ	413 غ	i 543	13غ
Faa	2305	87 ف	1213 ف	923 ن	82 ف
Gaaf	2784	97 ق	9 37 ق	1614 م	136 ق
Kaaf	2101	69 ك	5 914	988 ک	130 لك
Laam	6745	J 175	J 3546	1 2206	818 ل
Miim	7871	177 م	4 043 •	4 2844	r 807
Nuun	7484	2437 ن	; 1264	: 1905	1878 من
NuunChadda	225	0 نّ	3 0	1 225	0 ن
Haa	2670	۵ 223	a 704	4 1196	4 548
Waaw	4421	1 و	621	28 بو	300
Yaa	6641	317 ي	2 516	± 2640	ي 1168
YaaChadda	725	0 يّ	192	5 33	0 س
Hamza	192	ç 192			
HamzaAboveAlif	1437	i 1	102	i 3:	35
TaaaClosed	1417	441 ة			å 976
HamzaUnderAlif	253	1	182	Ļ7	1
AlifBroken	162	53 ی			109 س
TildAboveAlif	84	ĩ	32	Ĩ.5	2
HamzaAboveAlifBroken	210	3 ئ	i 167	i 39	1 سئ
HamzaAboveWaaw	89	30 ۇ		59 ـؤ	

Table 4: Distribution of letters in set 1

Letter label	Nb.	Isolate	Begin	Middle	End
٨lif	14925	15	1 5777		19
Baa	14925	150	. 2020	2244	220
Баа	4703	150 ب	2059	÷ 2344	÷ 244
Thee	622	042 ت 10	* 220	• 240	344 ב
	1907	19 ت 14	■ 230	± 349	دد ت در
Jiim	1897	ک 54	÷ /56	ب 1034	E 33
Haaa	2963	C 93	1159 ح	► 1619	で . ⁹²
Xaa	1435	č ¹⁸	622 خ	777 خ	Č ⁻¹⁸
Daal	3033	و د	963	20 د)70
Thaal	520	ا ذ	166	3 خ	54
Raa	6243	, 1	823	44 ر	20
Zaay	1054	; 3	379	6 ز	75
Siin	3556	77 س	ىد	عد	100 س
C1.''	1446	*	1338	2041	
Shiin	1446	22 ش	558 ىتد	838 عتد	28 س
Saad	1377	22 ص	420 صد	906 صد	29 ص
Daad	943	42 ض	374 ضد	492 ض	35 ض
Thaaa	1426	38 ط	401 ط	925 ط	62 ط
Таа	238	7 ظ	66 ظ	149 ط	16 ظ
Ayn	2823	E ⁸⁵	1074 ع	s 1543	2 121
Ghayn	970	ė 15	444 غ	è 495	ا غ 16
Faa	2256	62 ف	1184 ف	937 ن	73 ف
Gaaf	2734	104 ق	8 72 ق	1632 ھ	126 ق
Kaaf	2090	63 ك	891 ک	1002 ک	134 ك
Laam	6926	J 193	J 3513	1 2334	886 ل
Miim	7836	162 م	4 152 م	4 2704	8 18 •
Nuun	7433	2391 ن	3 1262	1 848	1932 من
NuunChadda	224	0 نّ	3 0	224	0 ن
Haa	2687	224 ه	a 705	4 1201	4 559
Waaw	4313	1 و	480	28 و	333
Yaa	6630	317 ي	2 432	± 2701	ي 1183
YaaChadda	727	0يّ	210	5 17	0 ق
Hamza	187	-	<u> </u>	187	ž
HamzaAboveAlif	1483	i 1	156	i 3	27
TaaaClosed	1407	429 ة			å 978
HamzaUnderAlif	250	[]	160	19	0
AlifBroken	161	, 4 7		÷	. 114
TildAboveAlif	84	ĩ	40	Ĩ 44	
HamzaAboveAlifBroken	208	0 ئ	i 166	i 34	8 مئ
HamzaAboveWaaw	90	ؤ	32	ـ ـ ـ ـ ـ ـ ـ ـ ـ ـ ـ ـ ـ ـ ـ ـ ـ ـ ـ	58

Table 5: Distribution of letters in set 2

Letter label	Nb. Occ	Isolate	Begin	Middle	End
Alif	15165	15	988	L 91	77
Baa	4692	156 ب	1955	÷2343	238 ب
Тааа	9897	617 ت	33546	: 5380	354 ت
Thaa	631	16 ث	3 245	3 35	35 ث
Jiim	1887	د 53	ج 784	> 998	52
Haaa	3017	ट 63	> 1194	> 1659	ح- 101
Xaa	1439	τ ¹¹	┶ 643	š 765	خ 20
Daal	3075	و د	947	21 د	128
Thaal	528	:	185	i 3	43
Raa	6169	1 ر	746	44 ر	23
Zaay	1054	; 3	362	; 6	92
Siin	3674	75 س	ىىد	عد	103 س
			1411	2085	
Shiin	1418	18 ش	545 شد	827 عثد	28 مش
Saad	1388	17 ص	390 صد	948 ص	33 ص
Daad	936	50 ض	346 ضد	511 ض	29 ض
Thaaa	1431	39 ط	393 ط	937 <u>ط</u>	62 ط
Таа	240	1 ظ	46 ظ	176 ظ	17 ظ
Ayn	2769	E 64	1015 ع	s 1560	E 130
Ghayn	983	<u>1</u> 2 غ	423 غ	غ 530	اع اغ
Faa	2221	54 ف	1178 ف	à 910	79 ف
Gaaf	2853	107 ق	984 ق	1640 م	122 ق
Kaaf	2099	76 ك	S 904	996 ک	123 ك
Laam	6972	j 183	J 3606	1 2259	924 ل
Miim	7957	190 م	4 066 4	4 2899	802
Nuun	7289	2319 ن	: 1293	: 1811	1866 من
NuunChadda	224	0 نّ	0 ڏ	i 224	0 ٽ
Haa	2590	192 ه	▲ 631	4 1222	4 546
Waaw	4325	1 و	507	28 و	818
Yaa	6876	318 ي	2 2527	± 2764	ي 1270
YaaChadda	709	0 يّ	198	4 511	0 س
Hamza	190		, , , 1		
HamzaAboveAlif	1455	11	133	i 3:	22
TaaaClosed	1394	435 ة			ä 959
HamzaUnderAlif	256	Į	169	ļ 8	37
AlifBroken	164	58 ی		-	106 س
TildAboveAlif	83	ĩ	39	ĩ 44	
HamzaAboveAlifBroken	208	4 ئ	170 ئ	1 27	7 مئ
HamzaAboveWaaw	89	ۇ	21	_ؤ	68

Table 6: Distribution of letters in set 3

Letter label	Nb.	Isolate	Begin	Middle	End
Alif	Occ 15120	15	866	1.92	54
Baa	4704	132	. 1070	. 2362	
Тааа	9797	÷ 633	3625	÷ 5202	÷ 331
Thaa	634	- 033	• 3023	▲ 3208 * 260	÷ 26
liim	1024	- 61	· 202	▲ 300 • 1016	- 20
Hana	2022	E 01	÷ 000	× 1010	æ 39
Пааа Vaa	2933	C 08	► 1205	× 1552	¹⁰⁸ ح
	1401	Č ¹⁶	► 615	▲ /49	Č -21
Daal	2990	د	7 09	1 20	181
Inaal	504	د ا	144	- 3	60
Raa	6335	ا ر	833	45 و	02
Zaay	1066	j 4	100	6 ز	66
Siin	3512	63 س	ىد		94 س
Shiin	1/3/	÷. 17	1349	2006	÷. 25
Sand	1411	11 دس	J90		22 س
Daad	006	19 ص	• 201	· 457	33 ص
Daau	900	34 ص	اەد صد	45/	34 <u>ص</u>
Thaaa	1420	34 طد م	399 ط	929 <u>ط</u>	64 <u>ط</u>
Taa	238	0 طد	64 ط	159 ط	15 ط
Ayn	2/18	E 72	1016 ع	a 1518	E ¹¹²
Ghayn	984	É 12	399 غ	± 566	٤ ⁷
Faa	2313	73 ف	1264 ف	894 ف	82 ف
Gaaf	2883	106 ق	999 ق	1639 ه	139 ق
Kaaf	2145	86 ك	935 ک	978 ک	146 ك
Laam	7002	J 207	J 3656	1 2247	892 ل
Miim	7806	157 م	a 3963	 2848	A 838
Nuun	7316	2341 ن	; 1239	i 1860	1876 من
NuunChadda	223	0 نّ	0 ڏ	223	0 بنّ
Наа	2718	201 ه	▲ 681	4 1252	4 585
Waaw	4333	1 و	494	28 و	339
Yaa	6685	322 ي	2443 ي	± 2699	ي 1221
YaaChadda	719	0 يّ	215 2	- 504	0 س
Hamza	193	<u>ج</u> اري		193	
HamzaAboveAlif	1512	i 1164 i 348		48	
TaaaClosed	1364	398 ٽ			å 966
HamzaUnderAlif	247	<u>į</u> 1	71	Ļ7	6
AlifBroken	163	ع 42 ي		-	121 س
TildAboveAlif	83	Ĩ	38	Ĩ.4	5
HamzaAboveAlifBroken	209	5 ئ	i 161	i 35	8 مئ
HamzaAboveWaaw	91	ۇ	24	_ؤ	57

Table 7: Distribution of letters in set 4

Letter label	Nb.	Isolate	Begin	Middle	End	Letter label	Nb.	Isolate	Begin	Middle
Alif	15046	15	689	19	357	Alif	15019	15	797	19
Baa	4730	. 161	1991	. 2341		Baa	4717	. 146	1998	. 2354
Гааа	9942	بة 101 ت 580	3629	: 5389	<u>.</u> 344	Тааа	9897	ب 641 ن	3612	÷ 200 1
Thaa	643	<u>ے</u> 26	3 242	1 347	<u>ے</u> 28	Thaa	628	- 011 ش 22	3 227	1 353
Jiim	1915	- 20	▲ 809	► 990	7 56	liim	1939	- 49	▲ 803	- 1048
Нааа	3000	C 00	► 1134	► 1680	– 103	Нааа	3000	C 72	► 1180	► 1655
Xaa	1403	℃ 05	► 611	▲ 1000	$\frac{103}{23}$	Xaa	1407	$\dot{\boldsymbol{c}}^{03}$	► 618	▲ 1055 ▲ 751
Daal	3028	C 15	901	- 134	2127	Daal	3086	<u> </u>	30	- 751
Thaal	516		150		257	Thaal	518		164	
Paa	6253	د	024	<u>د</u>	420	Paa	6267	1	0 <i>4</i>	<u> </u>
Raa 7222	1042	ار	206	4 ر	429		1045) 1	804 277	ب ر (
Laay	3620	3	300	ا ر	70	Siin	3602	2		ېز ا
51111	3029	99 س	1401	2001	ہ/ س	SIIII	2002	3 / س	سر 1350	2062
Shiin	1455	2.5 ش	1401 566 شد	838 مثد	26 يىش	Shiin	1458	26 ش	1339 	817 مثر
Saad	1371	14 ص	413 ص	896 <u>مد</u>	<u>48 ص</u>	Saad	1389	<u>1</u> 2 ص	415 صد	921
Daad	921	41 ض	369 ضد	470 ض	41 ض	Daad	920	43 ض	335 ض	503 ض
Thaaa	1446	11 <u>س</u>	412 ط	934 ط	11 <u>س</u> 67 طد	Thaaa	1462	<u>بور الم</u>	ь 428	937 ط
Таа	239	<u>- 55</u>	<u>- 112</u>	169 ظ	<u> </u>	Таа	241	4 ط	5 طل	<u>158</u>
Avn	2755	c 68	1017 ء	• 1552	• 118	Avn	2723	<u> </u>	007 م	1 1510
Ghavn	990	č 15	ـ 1017 422 غ	• 534	• 19	Ghavn	1004	č 15	ـ 1007 425غ	▲ 540
Faa	2339	73 ف	i 1257	<u>م</u> 920	49 في 19	Faa	2315	<u>د م</u>	i 1226	<u>م</u> 928
Gaaf	2762	a 103	a 959	a 1574	a 126	Gaaf	2803	a 99	a 974	a 1584
Kaaf	2136	<u>ط</u> 84	5 914	S 980	128 بل	Kaaf	2140	<u>در م</u>	5 913	S 1004
Laam	6790	1188	13433	1 2288	1 881	Laam	6724	1174	t 3466	1 2203
Miim	7797	o 175	• 4067	• 2732	a 823	Miim	7817	▲ 166	4 038	• 2203
Nuun	7400	. 2435	1273	: 1825	1867	Nuun	7264	· 2411	1231	: 1835
NuunChadda	224	ر با م	- 12/3 30	1025	<u>ت، انتقارات</u>	NuunChadda	223	<u>تبع ر</u>	3 0	1000
Наа	2705	178 ه	▲ 699	a 1297	4 531	Наа	2724	00 230 ه	a 695	4 1236
Waaw	4264	- 1,0	466	• · · · · · · · · ·	798	Waaw	4352	1 م	514	o 1250
Yaa	6648	9 I	J 2507	+ 2656		Yaa	6735	9 I	J 2535	+ 2652
YaaChadda	735	¹²⁰ مي ا	<u>168</u>	<u>-</u> 567	<u> </u>	YaaChadda	733	501 ي	- 199	<u> </u>
Hamza	192	ں ي	1 100	192	• سي	Hamza	188	ي	1,79	188
HamzaAboveAlif	1456	11	158	172	208	HamzaAboveAlif	1427	11	113	100
TaaaClosed	1400	11	138	6 .	270	TaaaClosed	1385	× 120	115	6.3
HamzaIInderAlif	2/18	• 455	171	4	4 9/0	HamzaUnderAlif	247	• 43U	170	
AlifBrokon	161	. 55	1/1	Ļ	11		161	12	1/9	<u>ب</u>
Tild A hove A lif	101	ددی ت	16	7	¹⁰⁰ س 27	Tild Above Alif	101	⁴⁵ ی	27	7
Humze Above AllfDreiter	200	1	40	L 6.20	5/	Homzo A hove A lift rol	00		51	6.20
	208	2 ئ	¥ 16/	- 32	/ جئ	HamzaAboveAllBroken	210	6 ئ	164	ii 39
namzaAboveWaaw	89	هٔ ۱	28	ــة ١	.61	HamzaAboveWaaw	90	لة ا	23	ة I

Table 8: Distribution of letters in set 5

Table 9: Distribution of letters in set 6

1 سئ 67 ـؤ

4 Evaluation Protocols

In this section, we propose the definition of a set of robust benchmarking protocols on top of the APTI database. Preliminary experiments with a baseline recognition system have helped in calibrating and validating these protocols.. From the obtained results, we believe that the large number of data available in APTI and the different source of variability (cf Section 2.5) make it well suited for significant and challenging evaluation of systems.

4.1 Error estimation

The objective of any benchmarking of recognition systems is to estimate, as reliably as possible, the classification error rate \hat{P}_e . It is important to keep in mind that, whatever the task and data used, \hat{P}_e is a function of the split of the data into training and test sets. Different splits will result in different error estimates. Hopefully, APTI is composed of quite large sets of data, which is helping in reaching stable estimates of \hat{P}_e .

Our objective is then to obtain a reliable estimate of \hat{P}_e while keeping the computation load tractable. Therefore, we have opted for a *rotation method*, as described in [Jain 00, Section 7]. The idea is to reach a trade-off between the *holdout method* which leads to pessimistic and biased values of the error rate and the *leave-one-out method* that gives a better estimate but at the cost of larger computational requirements. The rotation method we are proposing is illustrated in Fig. 6. The procedure is to perform independent runs on 5 different partitions between training and testing data.



Fig. 6: Illustration of the rotation method. For a given partition, the training sets are depicted in dark grey and the testing sets in light grey.

The final error estimate is taken as the average of the error rates obtained on the different partitions.

$$\hat{P}_{e} = \frac{1}{5} \sum_{i=1}^{5} \hat{P}_{e,i}$$

In the previous formula, $\hat{P}_{e,i}$ is the error rate obtained independently on a system trained and tested using the sets defined in partition *i*. The procedure actually corresponds to computing the average of performance of 5 independent systems.

4.2 Train and test conditions

Using the procedure described in section 4.1, we can define different combinations of train and test conditions. The objectives are to measure the impact of some of the variability of the data. We therefore propose 20 protocols as summarized in Table 3.

The notations Tr(font, style, size) and Te(font, style, size) define the training and testing conditions with:

- 1. the font label as indicated in Fig. 1
- 2. the style where *p*, *i*, *b* and *bi* are for plain, italic, bold and bold+italic
- 3. the size in points

We suggest researchers willing to define new protocols to use this notation to specify the conditions of their training and testing.

Protocol name	Train choice Tr(font, Style, Size)	Test choice Te (font, Style, Size)
APTI 1	Tr(B, p, 10)	Te(B, p, 10)
APTI 2	Tr(B, p, 10)	Te(B, i, 10)
APTI 3	Tr(B, p, 10)	Te(B, b, 10)
APTI 4	Tr(B, p, 10)	Te(B, bi, 10)
APTI 5	Tr(B, p, [6, 10, 14, 18])	Te(B, p, [6, 10, 14, 18])
APTI 6	Tr(B,[p,i,b], [6, 10, 14, 18])	Te(B,[p,i,b], [6, 10, 14, 18])
APTI 7	Tr([A,B,C,F,H], p, 10)	Te([A,B,C,F,H], p, 10)
APTI 8	Tr([D,E,G,I,J], p, 10)	Te([D,E,G,I,J], p, 10)
APTI 9	Tr([A,B,C,F,H], [p,i,b], 10)	Te([A,B,C,F,H], [p,i,b], 10)
APTI 10	Tr([D,E,G,I,J], [p,i,b], 10)	Te([D,E,G,I,J], [p,i,b], 10)
APTI 11	Tr([A,B,C], p, 10)	Te([F,H], p, 10)
APTI 12	Tr([D,E,G], p, 10)	Te([I,J],i, 10)
APTI 13	Tr([A,B,C], p,[6,10,14,18])	Te([F,H], p, [6,10,14,18])
APTI 14	Tr([D,E,G], p,[6,10,14,18])	Te([I,J], p, [6,10,14,18])
APTI 15	Tr(B, p, 6)	Te(B, p, 6)
APTI 16	Tr(B, p, 8)	Te(B, p, 8)
APTI 17	Tr(B, p, 10)	Te(B, p, 6)
APTI 18	Tr(B, p, 6)	Te(B, p, 10)
APTI 19	Tr(B, p, [6, 10, 14, 18])	Te(B, p, [7,9,12,24])
APTI 20	Tr(all, all, all)	Te(all, all, all)

 Table 3: APTI protocols

The objectives behind the protocols of Table 3 can be explained as follows:

- **APTI 1**: This is the baseline protocol where performances should be the highest as there are no mismatched between training and testing conditions.
- **APTI 2,3,4**: We measure here the capability of systems trained using plain style to generalize on italic, bold and bold+italic.
- **APTI 5,6**: While using the same font, we measure the capability of the system to treat different sizes.
- **APTI 7,8,9,10**: These experiments measure the capability of systems to recognize muti-font text.
- **APTI 11,12,13,14**: We measure the capability of systems to recognize unseen fonts text.
- **APTI 1,15,16,17,18,19**: Firstly, we measure the potential degradation of performance using smaller sizes. Secondly, we measure the capability to recognize unseen sizes.
- **APTI 20**: This is the global experiment where all available data is used for training and testing.

5 Conclusions

APTI, a new large Arabic printed text images database is presented together with evaluation protocols. APTI aims at the large-scale benchmarking of open-vocabulary text recognition systems. While it can be used for the evaluation of any OCR systems, APTI is, by nature, well suited for the evaluation of screen-based OCR systems. The challenges addressed by the database are in the variability of the sizes, fonts and style and the protocols that are defined are crafted to put into evidence the impact of such variability. APTI will be made publicly available for the purpose of research.

6 References

[Pechwitz 02]	M. Pechwitz, S. S. Maddouri, V. Maergner, N. Ellouze, and H. Amiri. IFN/ENIT - database of handwritten Arabic words. In Proc. of CIFED 2002, pages 129–136, Hammamet, Tunisia, October 21-23 2002
[Slimane 08]	F. Slimane, R. Ingold, M. A. Alimi and J. Hennebert, Duration Models for Arabic Text Recognition using Hidden Markov Models. CIMCA 2008, Vienne, Austria, December 10-12 2008
[Jain 00]	A. K. Jain, R. Duin and J. Mao, Statistical Pattern Recognition: A Review, IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 22, No. 1, January 2000
[Khorsheed 07]	M. S. Khorsheed, Offline recognition of omnifont Arabic text using the HMM ToolKit (HTK). Pattern Recognition Letters 28(12): 1563-1571, 2007
[Schlosser 95]	S. Schlosser, "ERIM Arabic Database", Document Processing Research Program, Information and Materials Applications Laboratory, Environmental Research Institute of Michigan, 1995
[Margner 05]	V. Margner, M. Pechwitz, H. El Abed, "Arabic Handwriting Recognition Competition", In ICDAR, 2005, pp.70 – 74
[Margner 07]	V. Margner and H. E. Abed. "ICDAR 2007 Arabic handwriting recognition competition". In ICDAR, Sept. 2007 vol. 2, pp. 1274–1278.
[Graff 06]	D. Graff, K. Chen, J. Kong, and K. Maeda, "Arabic Gigaword Second Edition", Linguistic Data Consortium, Philadelphia, 2006
[Abbes 04]	R. Abbes, J.D. Hassoun, "The Architecture of a Standard Arabic Lexical Database, Some Figures, Ratios and Categories from the DIINAR.1 Source Program", Workshop of Computational Approaches to Arabic Script-Based Languages, Geneva, 2004
[Husni 08]	Husni A. Al-Muhtaseb, Sabri A. Mahmoud, Rami S. Qahwaji, Recognition of off-line printed Arabic text using Hidden Markov Models. European Signal Processing Conference. Vol. 88, Issue 12, Pages 2902-2912, Lausanne, Switzerland, August 25-29, 2008

[Shaaban 08]	Z. Shaaban, A New Recognition Scheme for Machine-Printed Arabic Texts based on Neural Networks. Proceedings of World Academy of Science, Engineering and Technology, Vol. 31, Vienna, Austria, July 25-27 2008
[AbdelRaouf 08]	A. AbdelRaouf, C. A Higgins, and M. Khalil, A Database for Arabic Printed Character Recognition. ICIAR 2008, LNCS 5112, pages 567–578, 2008.
[Kanoun 2005]	S. Kanoun, A. M. Alimi, Y. Lecourtier, "Affixal approach for Arabic decomposable vocabulary recognition a validation on printed word in only one font", In ICDAR, Sept. 2005, vol. 2 pp.1025 - 1029
[Baird 08]	H. S. Baird. "State of the Art of Document Image Degradation Modeling". Proceedings of the 4th IAPR Workshop on Document Analysis Systems, DAS 2000.