

Database Specification

Overview

- Over 2200 form-pages of (937) Tunisian town/village names
- All data are digitized using a resolution of 300 dpi (b/w)
- Ground truth information available
 - e.g. sequence of Arabic character shapes
 - baseline/reference line position
 - topline position for set_a
- A wide variety of writing styles; 411 different writers
- About 26000 Arabic handwritten Tunisian town/village names
- Approximately 212 000 Arabic characters and ligatures
- Divided into 4 sets (a-d)
- Images and ground truth documentation included

Set specification

The whole database is divided into 4 disjoint sets for training and testing Arabic OCR systems. We recommend using the specified sets for comparability with results of other groups.

SET	Number of words	Number of writer (OK+bad)
a	6537	88+14=102
b	6710	89+13=102
c	6477	88+15=103
d	6735	90+14=104
SUM	26459	355+56= 411

Data formats

Image format

The form-pages are stored in uncompressed TIFF-file format (file-extension .tif).

All cropped Tunisian town/village names coming as uncompressed TIFF-images and as BMP-images (file-extension .bmp). In the TIFF header section “Image description” label information is stored.

Ground Truth format

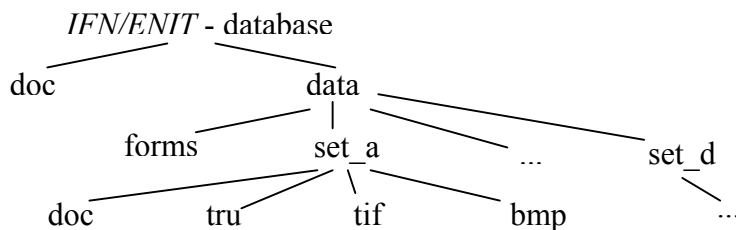
For each cropped Tunisian town/village name you find a truth file (file-extension .tru). The truth-file is an ASCII .txt file including all available ground truth information. One example is shown below.

```
01: COM: IFN/ENIT-database truth (label) file
02: COM: http://www.ifnenit.com
03: COM: IfN, TU-BS
04: COM: di45_019.tif coming from pb377_6.tif
05: X_Y: 498 87
06: BDR: begin data record
07: LBL: ZIP:3032;AW1:مركز درويش;AW2:maB|raE|keB|zaE|daA|raA|waA|yaB|shE|;QUA:YB1;ADD:P6
08: CHA: 9
09: BLN: 56,42
10: TLN: 23,19
11: EDR: end of data record
```

lines 01-04: comments
line 05: image size in pixel (x,y)
line 06: begin data record
line 07: label
ZIP: Tunisian post code / ZIP code
AW1: Tunisian town/village name in Arabic windows encoding
AW2: Arabic character shape sequence of Tunisian town/village name in Latin code (refer Appendix for the lookup table).
QUA: baseline quality tag (B1=:OK;B2=:bad)
ADD: number of pieces of Arabic words (PAW's)
line 08: number of characters
line 09: baseline/reference line information Y1,Y2
line 10: topline information Y1,Y2 (in data set_a only!!!)
line 11: end of data record

Database Organisation

The IFN/ENIT-database has the following directory structure.



doc

In this directory documentation in pdf and/or txt format is available.

data

In this directory all available data are stored. Each of the four sets has it's own directory (set_?).

set_?

In these directories all data are available in tif and bmp file format. The ground truth information you will find in the directory tru. The subdirectory doc under data/set_? includes a pdf-file for each writer with all words and baselines.

File name convention

The following file name convention is used.

SW_{ww}_NNN.EXT

- S:=set (a,b,c,d)
- W_{ww}:=writerID; W=(e,f,i,j,m,q) ,w=(0..9)
- NNN:=word_number; N=(0..9)

Ordering Information

IFN/ENIT-database is made available for non-commercial use. The data is supplied with no guarantee of accuracy or usability. We can't guarantee to maintain the IFN/ENIT-database, but would be interested in hearing of any comments or results that you have.

Upon request we make the data available on the Internet for free download. If the database has to be shipped as a CD-Rom production and shipping costs will be charged. In both cases please contact us.

BUGS

When you are working with the *IFN/ENIT*-database perhaps you will discover bugs concerning the label or the extraction of the data. Please report the bugs you find back to us. So we can improve the quality of the data over the time.

Reporting results

We kindly invite you to publish reached results with the *IFN/ENIT*-database. To keep recognition results comparable, we suggest reporting results as shown in the following example:

Database version: <i>IFN/ENIT</i> -database v1.0p2			
Test	Training set(s)	Test set	Recognition result(*)
1	a,b,c	D	23.7%
2	c,b,d	A	25,4%
...			

(*) Percentage of correctly recognised words in the specified test set. We recommend to use the ground truth / label category “ZIP:” as reference for the recognition result. Please use for each test the whole number of 937 different Tunisian town/village names as lexicon, in the case a lexicon is needed.

Contact

Please feel free to contact contact@ifnenit.com.

References

Mario Pechwitz, Samia Snoussi Maddouri, Volker Märgner, Noureddine Ellouze, Hamid Amiri; *IFN/ENIT-database of handwritten Arabic words*, In Proceedings of CIFED'02, Hammamet, Tunisia, 21.-23.10.2002, p. 129-136

Appendix

Label legend / lookup table – Arabic label to Latin label & statistic of occurrence in the database

Arabic label	Latin label	Quantity
0 A	0A	342
1 A	1A	279
2 A	2A	384
6 A	6A	311
7 A	7A	354
8 A	8A	284
9 A	9A	341
هـ A	hhA	520
آ A	amA	544
إ A	aeA	1660
أ A	ahA	631
أ E J B	aeElaB	360
أ E J B	ahElaB	122
م ع	alM	355

أ A	aaA	20251
أ E	aaE	13308
أ E J B	aaElaB	799
أ E J M	aaElaM	1076
ب A	baA	331
ب B	baB	5636
ب E	baE	344
ب M	baM	3407
ة A	teA	2182
ت A	taA	356
ت B	taB	2324
ة E	teE	7259
ت E	taE	357
ت M	taM	1045

ث A	thA	338
ث B	thB	353
ث M	thM	327
ج A	jaA	504
ج B	jaB	981
ج E	jaE	346
ج M	jaM	1218
ج MJ B	jaMlaB	539
ح A	haA	314
ح B	haB	2483
ح E	haE	295
ح M	haM	1804
ح MJ B	haMlaB	365
ح M _م MJ B	haMmaMlaB	64
ح M _ن B	haMnaB	100
خ A	khA	341
خ B	khB	863
خ E	khE	321
خ M	khM	425
خ MJ B	khMlaB	310
د A	daA	2718
د E	daE	4883
ذ A	dhA	353
ذ E	dhE	703
ر A	raA	6252
ر E	raE	9369
ز A	zaA	1764
ز E	zaE	2718
س A	seA	873
س B	seB	4218
س E	seE	811
س M	seM	1110
ش A	shA	475
ش B	shB	1617
ش E	shE	351
ش M	shM	1277
ص A	saA	355
ص B	saB	956
ص E	saE	357
ص M	saM	1096
ض A	deA	351
ض B	deB	730
ض E	deE	343
ض M	deM	328

ط A	toA	343
ط B	toB	359
ط E	toE	350
ط M	toM	1258
ظ B	zaB	339
ظ M	zaM	690
ع A	ayA	915
ع B	ayB	1650
ع E	ayE	347
ع M	ayM	1990
غ B	ghB	326
غ M	ghM	600
ف A	faA	319
ف B	faB	898
ف E	faE	316
ف M	faM	1647
ق A	kaA	397
ق B	kaB	2608
ق E	kaE	348
ق M	kaM	1307
ك B	keB	1221
ك E	keE	335
ك M	keM	980
ل A	laA	1485
ل B	laB	14340
ل E	laE	1056
ل M	laM	2594
م A	maA	890
م B	maB	3886
م E	maE	536
م M	maM	4626
م MJ B	maMlaB	458
ن A	naA	1267
ن B	naB	3723
ن E	naE	2119
ن M	naM	2912
ه A	heA	696
ه B	heB	1924
ه E	heE	351
ه M	heM	347
و A	waA	3511
و E	waE	6529
ى A	eeA	322
ي A	yaA	2932

ي B	yaB	4383
ي E	eeE	350
ي E	yaE	2167
ي M	yaM	7759

Note:

- “llL” is often added and means ligature “chadda”
- The character shape indicators (A,B,M,E) are sometimes supplemented with a “1” or a “2”, like ”baA1”. In this case there is a point error detected. Take care: this feature is not consistent over the whole database ☹. We recommend ignoring these kinds of label supplements.