

# URDU Optical Character Recognizer

Benazir Mumtaz, Hina Khalid, Muhammad Kamran

ISSALE 2014

September 26, 2014

- Pre-processing
  - Line segmentation
  - Ligature segmentation
  - Connected Component identification
- Recognition
- Post-processing
  - Ligature formation

پر اسے کے ابا بابا پایا نایاب اب بابا ہے بات پایا پار کا کے  
 میں اب کے کا کا راز کا ہے راد وا کے دیا وار ہے داد ہے  
 میں کے سے راد وا سے دیا رب زار سے پر برا کے فڈ ہے  
 میں کے کا کا پڑ بڑا پر اب کے رکر کر کر ہے وہ ہے کا  
 میں بات کے ہے سے رو کر ڈرہار ہے کا کا سے کر میں ہے  
 میں کے کی کی کا کی رکر باب سے میں داد میں ں کان بڑ پر  
 رں ناں کاں پرن برن پڑن بڑر زار زر سے کرنا سے  
 سے اور سے کے و کے کے کے کے کے کے



- ligature Shapes : 15 main strokes

ا ب دے کا کر کے کی مس ر و ر سے ن  
300 301 302 303 304 305 306 307 308 109 310 311 312 313 314

- diacritics : 5 diacritics
- lines: 8 lines

Main Body Type	Total Tokens in Document Image	Total Unique Syllable in Document Image
300	18	1
301	18	5
302	8	2
303	9	3
304	11	1
305	12	1
306	8	1
307	10	1
308	8	1
309	8	1
310	12	4
311	11	1
312	23	2
313	9	1

- Line Segmentation: [8:8] = 100%
- Ligature Segmentation: [179:180]=100%
- Connected Component Disambiguation:

MB in Document Image	MB Detected	Diacritics in Document Image	Diacritic Detected
180	180	88	88

Training Accuracy : 97.5 %

```

BB_Area <= 1020
  width <= 17
    width <= 9: three_hundred (35.0)
    width > 9
      Height <= 31: three_hundred_eleven (35.0)
      Height > 31: three_hundred_one (35.0)
  width > 17
    width <= 21
      Height <= 25: three_hundred_twelve (35.0)
      Height > 25: three_hundred_three (35.0)
    width > 21
      Height <= 22: three_hundred_two (35.0)
      Height > 22: three_hundred_ten (35.0)
BB_Area > 1020
  Height <= 45
    width <= 58: three_hundred_four (56.0)
    width > 58
      Height <= 27: three_hundred_thirteen (28.0)
      Height > 27
        BPC <= 530: three_hundred_four (11.0)
        BPC > 530
          Height <= 29: three_hundred_thirteen
(5.0)
    
```



Training Accuracy : 91.8 %

Height <= 14

width <= 14: three\_hundred\_twenty (36.0/1.0)

width > 14: three\_hundred\_twentyone (35.0)

Height > 14

width <= 14

Height <= 16: three\_hundred\_twentyfive (19.0)

Height > 16

Height <= 17

BB\_Area <= 209

BPC <= 77: three\_hundred\_twentyfour

(10.0)

BPC > 77: three\_hundred\_twentyfive

(5.0/2.0)

BB\_Area > 209: three\_hundred\_twentyfive

(18.0/5.0)

Height > 17: three\_hundred\_twentyfour (18.0)

width > 14: three\_hundred\_twentythree (19.0)

Class Type	Total Samples	Correctly Recognized	% Accuracy
301	15	3	20%
302	15	2	13%
303	15	3	20%
305	15	15	100%
306	15	15	100%
307	15	15	100%
308	15	15	100%
309	15	15	100%
310	15	15	100%
311	15	0	0
312	15	3	20%
314	15	15	100%

Class Type	Total Samples	Correctly Recognized	% Accuracy
300	18	18	100
301	18	0	0
302	8	8	100
303	9	8	88
304	11	8	72
305	11	11	100
306	8	6	75
307	10	10	100
308	8	8	100
309	8	8	100
310	12	2	16
311	11	6	55
312	23	18	78
313	9	9	100
314	8	4	50

**310**

**320**

ب

**303**

**320**

ذ

**303**

**325**

ڈ

**325**

**310**

**323**

پ

**325**

**310**

**320**

پ

**306**

ک

**301**

**324**

ہا

**308**

کی

## Tesseract vs Decision Trees

64.45% — 63.75 %