

Muna Khayyat (PhD Student – 3<sup>rd</sup> Year – Concordia University )  
Center for Pattern Recognition and Machine Learning (CENPARMI)  
Supervisors: Dr. C. Y. Suen and Dr. L. Lam

## Introduction

Word spotting has been widely implemented for Latin-based and Chinese documents. However, few word spotting systems have been implemented for Arabic handwritten documents. Yet, Arabic is spoken by a significant number of the world's population. Arabic script is cursive by nature; besides, in Arabic writing words have no clear boundaries; these facts make the implementation of word spotting for Arabic handwritten documents a significant challenge.

We proposed a learning-based word spotting system that uses Support Vector Machine (SVM) to recognize sub-words rather than complete words. We will use this partial segmentation concept to resolve the word boundary problem in Arabic handwriting. In addition, a shortest distance algorithm is used to spot Arabic handwritten words.

## Work Progress

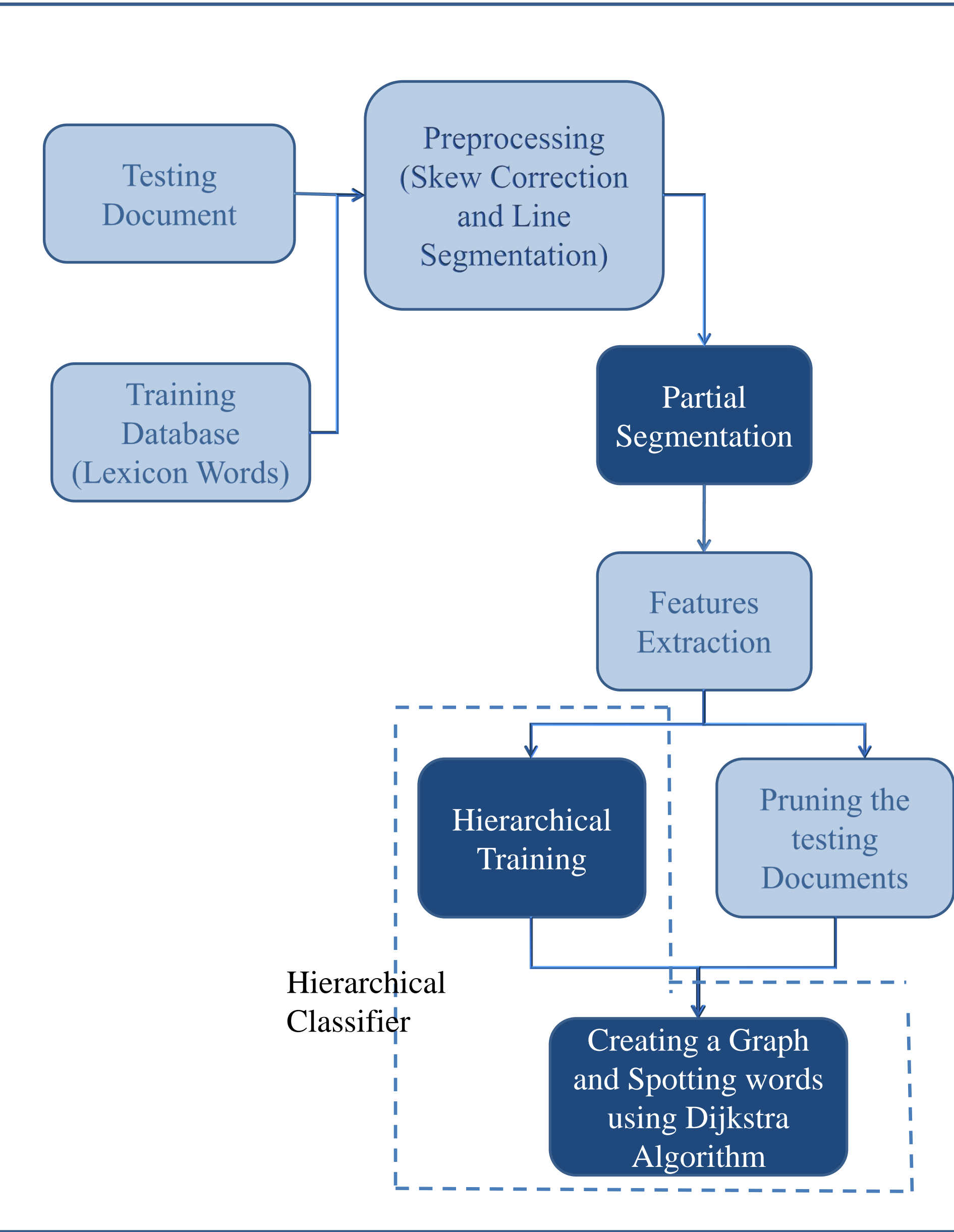
### Progress to date

1. An Arabic word recognition system using three different feature sets has been constructed. The features were tested on CENPARMI and other databases. The highest classification accuracy on CENPARMI Arabic words database is 95.46 % using gradient features. Another feature sets will be extracted in the future. Those features are supposed to bring higher classification accuracy.
2. Partial segmentation of the words to their PAWs.
3. Arabic text line segmentation.

### Current Work

Document Pruning.

## Research Plan

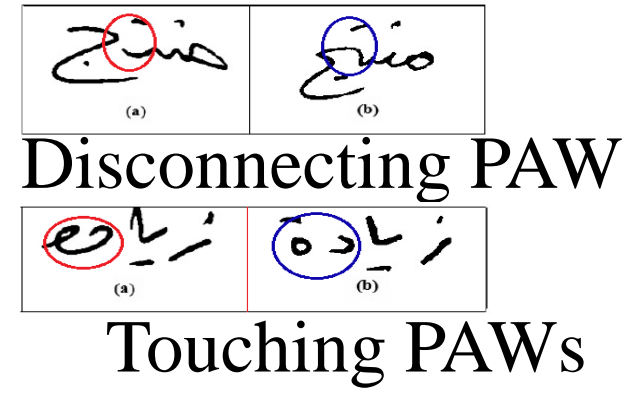


## Partial Segmentation

We reconstructed CENPARMI database [1] by segmenting it into PAWs. Some PAWs may be touched or disconnected . However, people seems to connected or disconnect almost in a similar manner. Thus, additional classes were considered for the new formulated PAW. The machine will learn these new classes instead of applying some heuristics and rules trying to connect the disconnected PAW and vise versa.

### Results

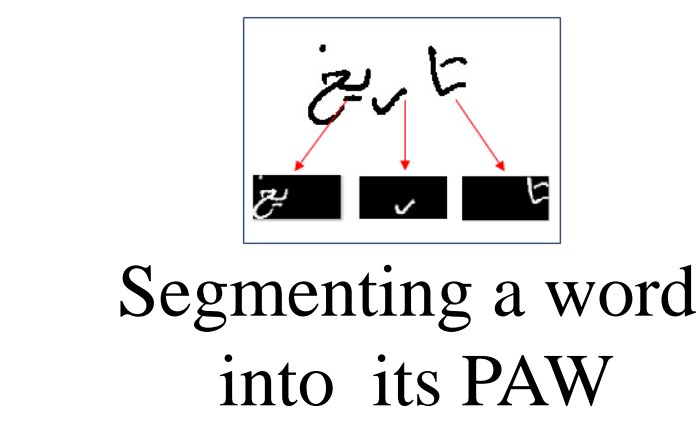
#### Difficulties



| Class Number | PAW | Confusing class | PAW |
|--------------|-----|-----------------|-----|
| 3            | ر   | 15              | ا   |
| 3            | ر   | 19              | هـ  |
| 3            | ر   | 23              | ز   |
| 10           | ن   | 45              | ن   |
| 17           | ن   | 44              | س   |
| 33           | س   | 34              | ب   |
| 58           | س   | 73              | س   |

Confusing PAW in the database

#### Partial Segmentation



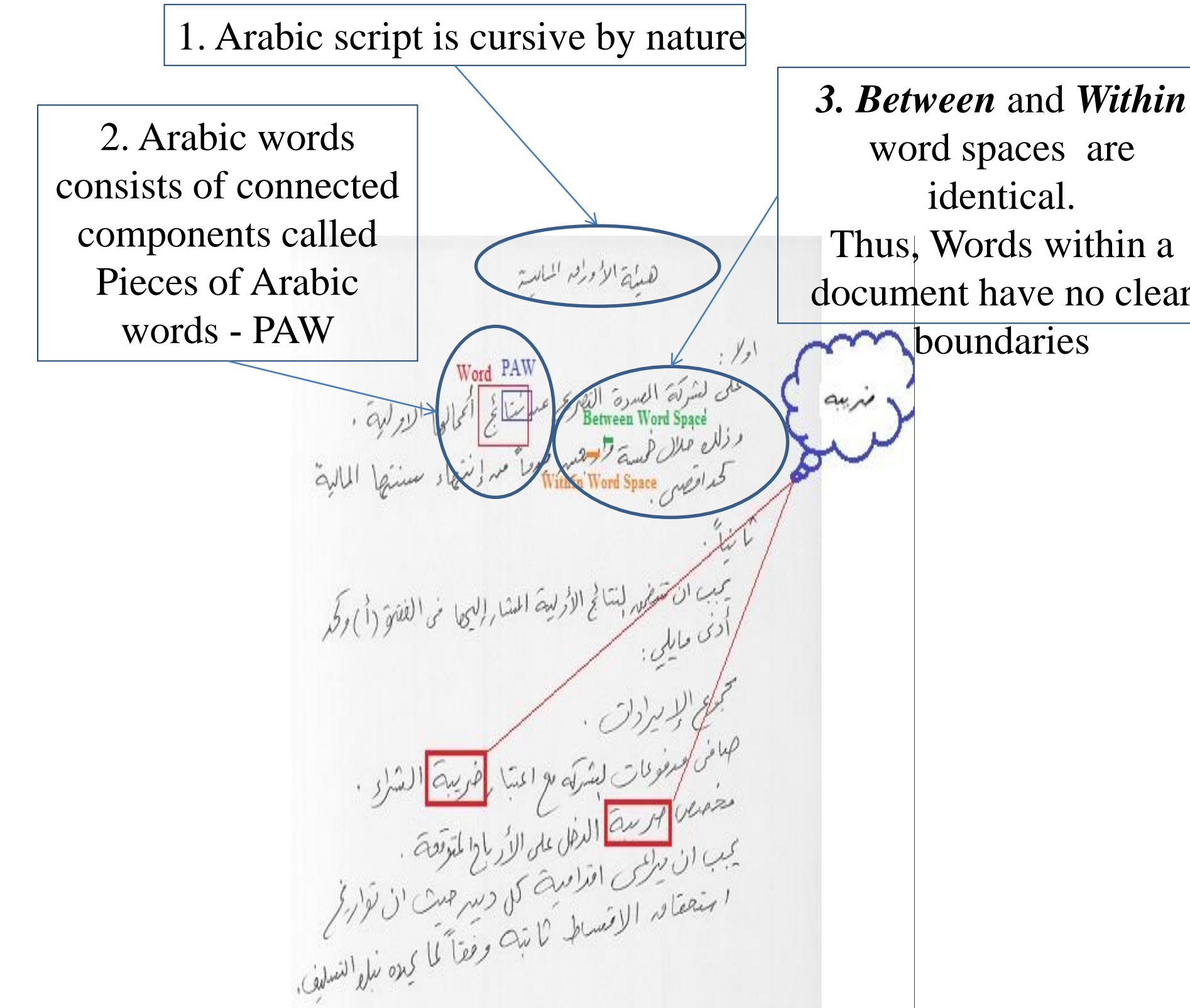
Segmenting a word into its PAW

#### PAW Recognition Results

| Class ID | Printed Sample | Confusing Class ID | Printed Sample of the confusing class | Number of misclassified samples | Misclassified samples in the database |
|----------|----------------|--------------------|---------------------------------------|---------------------------------|---------------------------------------|
| 60       | ت              | 10                 | ن                                     | 28                              |                                       |
| 3        | ر              | 19                 | هـ                                    | 20                              |                                       |
| 19       | هـ             | 26                 | و                                     | 19                              |                                       |
| 10       | ن              | 45                 | ن                                     | 18                              |                                       |
| 19       | هـ             | 3                  | ر                                     | 17                              |                                       |
| 26       | و              | 19                 | هـ                                    | 16                              |                                       |
| 25       | ف              | 10                 | ن                                     | 15                              |                                       |
| 23       | ز              | 7                  | م                                     | 11                              |                                       |
| 82       | س              | 42                 | س                                     | 10                              |                                       |

The most frequent confusion classes taken from the confusion matrix resulted from extracting the gradient features from the PAW images.

## Problems with Arabic Word Spotting



The aforementioned facts makes it difficult to segment Arabic handwritten document into words

We have to find a way to segment the document into PAWs and then reconstruct the words

**Hierarchical classifier**

## Arabic Handwritten Word Recognition

Three sets of features were extracted to recognize Arabic handwritten words, with each set of feature passed to one classifier. (Support Vector Machine - SVM) The confidence levels and classification results of the classifiers were used for the final classification. The CENPARMI database is used for validation.

### Preprocessing

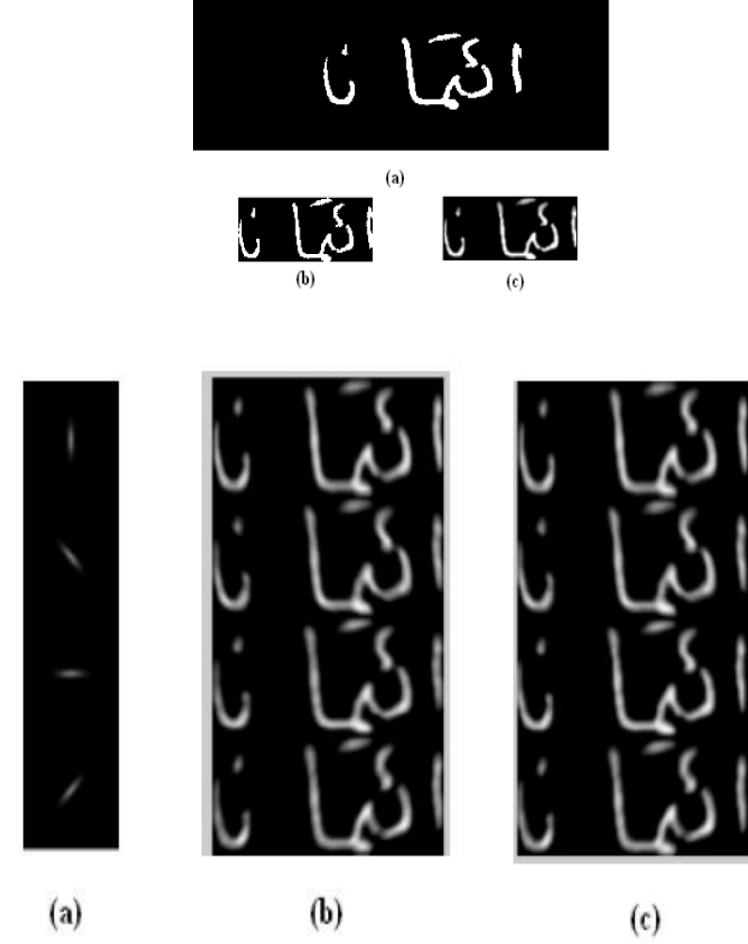
1. Binarization (Otsu Algorithm)
2. Size normalization  $120 \times 50$
3. Smoothed grayscale Image

### Feature Extraction:

#### 1. Gradient Features [3]

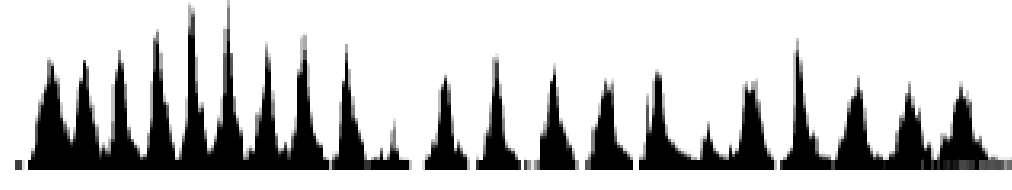
#### 2. Gabor Features [2]

#### 3. Fourier Features [4] The Discrete Fourier Transform of the projection, upper and lower profiles of the words were extracted:

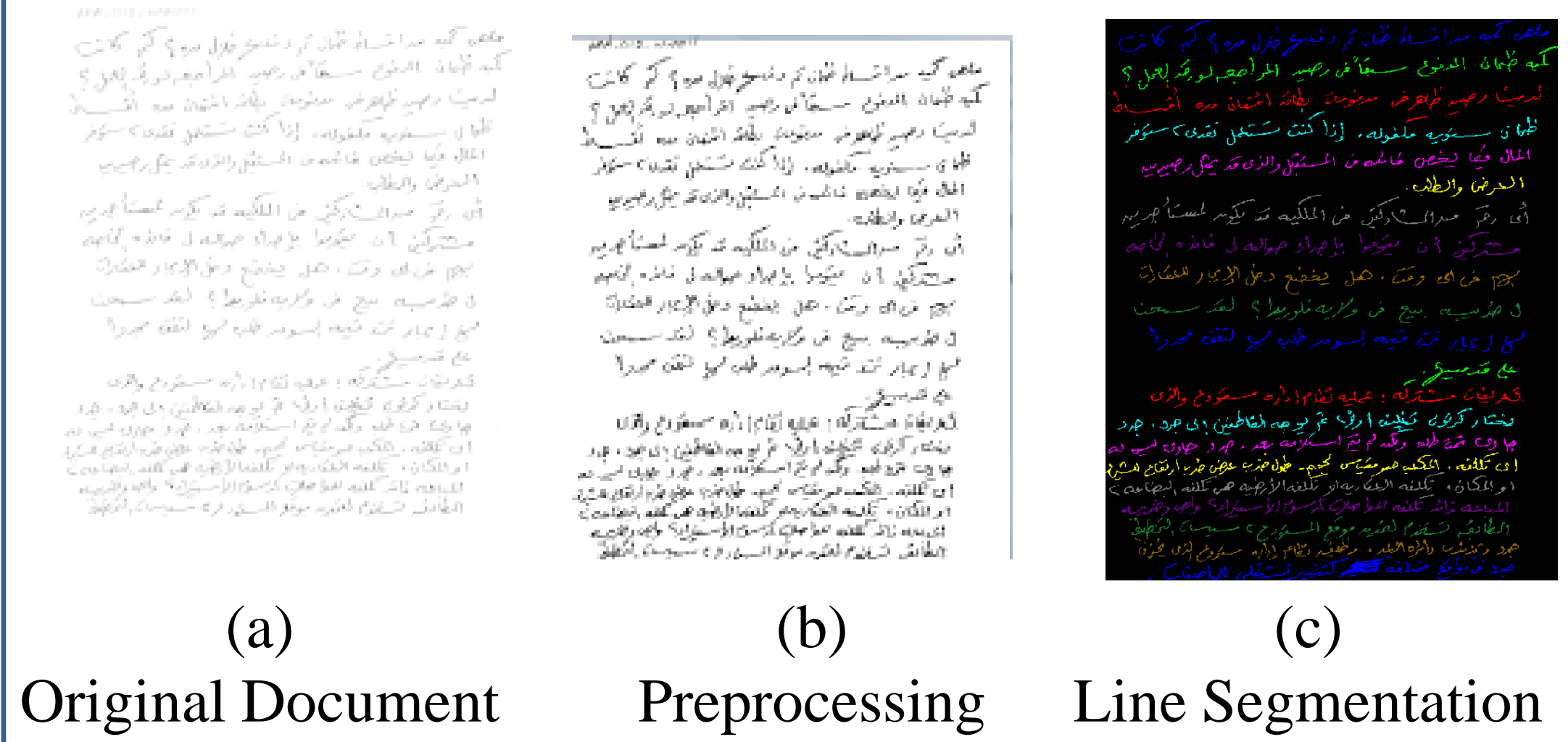


## Arabic Text Line Segmentation

Arabic handwriting is cursive by nature and it includes small connected components called diacritics. These facts add constrains to the segmentation of the Arabic handwritten documents into lines.



Probability Density Function (PDF) of the document



## References

- [1] H. Alamri, J. Sadri, C. Y. Suen, and N. Nobile “A Novel Comprehensive Database for Arabic Off-Line Handwriting Recognition,” In *Proceedings of Eleventh International Conference on Frontiers in Handwriting Recognition (ICFHR)*, pp. 664-669, Montreal, Canada, August 2008.
- [2] J. Chen, H. Cao, R. Prasad, A. Bhardwaj, and P. Natarajan, “Gabor Features for Online Arabic Handwriting Recognition”, In *Proceedings of the 9th International WORKSHOP on Document Analysis Systems (DAS '10)*, pp. 53 - 58, 2010.
- [3] M. Shi, Y. Fujisawa, T. Wakabayashi, and F. Kimura, “Handwritten Numeral Recognition Using Gradient and Curvature of Gray Scale Image”, *Pattern Recognition*, vol. 35, No. 10, pp. 2051 - 2059, 2002.
- [4] V. Lavrenko, T. M. Rath, and R. Manmatha, “Holistic Word Recognition for Handwritten Historical Documents”, In *Proceedings of the First International Workshop on Document Image Analysis for Libraries (DIAL)*, pp. 278 - 287, Palo Alto, CA, USA, January 2004.