

An Evaluation of Methods for Arabic Character Recognition

A. Lawgali

University of Benghazi, Libya
ahmed.lawgali@uob.edu.ly

Abstract

Off-line recognition of text plays a significant role in several applications such as the automatic sorting of postal mail or editing old documents. The recognition of Arabic handwriting characters is a difficult task owing to the similar appearance of some different characters. Most researchers have presented methods that recognise isolated characters. However, recognition of all shapes of Arabic handwritten characters still remains a great challenge. The selection of the methods for feature extraction and classification remain the most important step in achieving high recognition accuracy. The purpose of this paper is to compare the effectiveness of DCT and DWT in capturing discriminative features of all shapes of Arabic handwritten characters including overlapping characters with ANN and HMM in the classification stage. Since, the recognition of handwritten characters is an important step in the recognition of a word after segmentation, this paper ascertains the effectiveness of these techniques in capturing useful information and, hence, achieving more accurate recognition results. This work has been tested with HACDB database containing 6,600 shapes of Arabic characters. The results have demonstrated that the feature extraction by DCT with ANN yields a higher recognition rate.

Keywords: Arabic character, DCT, DWT, ANN, HMM

1. Introduction

Automatic off-line Arabic handwriting recognition still faces great challenges. Little research has been carried out in the field of Arabic handwritten character recognition compared to research into its Latin and Chinese counterparts. Arabic script is written from right to left and is composed of 28 characters, with no upper or lower case. Each character has two or four shapes, the shape of the character depends on its position in the word. Many languages use Arabic characters such as Persian, Urdu and Jawi [1]. The shape of some characters is similar but the difference arises with the position and the number of dots. Some handwritten characters may appear to be similar although they are different and it is difficult for the human eye to spot the difference [2]. In an automatic recognition system, the selection of the methods for feature extraction and classification might be the most important step for achieving a high recognition accuracy. Alma'adeed *et al.* [3] presented a system for recognition of handwritten Arabic words based on Hidden Markov Model (HMM). El-Hajj *et al.* [4] proposed a system using baseline dependant features and HMM for classifying handwritten Arabic words. Alma'adeed [5] introduced a system using Artificial Neural Network (ANN) for unconstrained handwritten Arabic words. Alkhateeb *et al.* [6] presented a technique for the recognition of handwritten Arabic words where Discrete Cosine Transform (DCT) is used for extracting features of the word. These features are then fed into a neural network for classification. Common methods in the classification stage are the use of an

ANN and HMMs [7]. A comparison between DCT and DWT to capture features of Arabic handwritten characters without overlapping characters has been introduced by Lawgali *et al.*[8]. Both DCT and DWT are widely used in the field of digital signal processing applications [9]. Therefore, in this paper, ANN and HMM are used to compare their effectiveness in classifying shapes of handwritten Arabic characters including overlapping characters via the features captured by DCT and DWT. This is an important step for the recognition of processes after segmentation.

The organization of the paper is as follows. Section 2 describes data acquisition and pre-processing. Section 3 discusses the two methods used for feature extraction and the classification stage is described in section 4. Section 5 discusses the results and their analysis. Section 6 concludes the paper.

2. Data Acquisition and Pre-processing

Arabic characters were read from HACDB database [10] which contains 6,600 shapes of handwritten Arabic characters, including overlapping characters as one shape. A form for collecting data was designed with 66 small squares divided into 6×11 small squares. Each writer filled in two forms with all the shapes of Arabic characters: as isolated, at the beginning of the word, in the middle, and at the end of the word; also, the shapes of the overlapping characters but without dots. Figure [1] shows a form filled in by one writer. The pre-processing task is used to remove the details that have no discriminative power in the process of recognition (*i.e.*, redundant data). Noise removal, binarization and normalization were carried out in the development of the database [10].

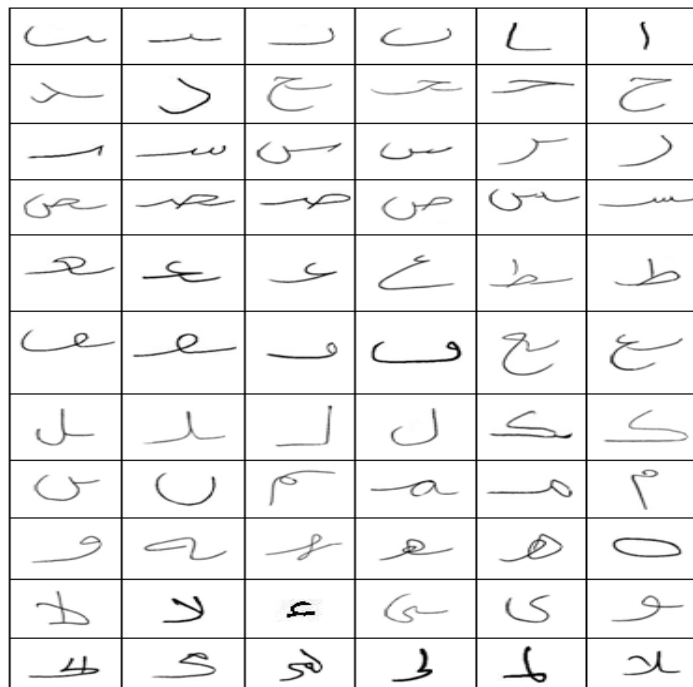


Figure 1. A form designed by the author to collect shapes of handwritten Arabic characters

3. Feature extraction

Feature extraction converts the image into a set of vectors to be passed onto the matcher to assist in the classification process. Therefore, these features should possess the essential characteristics of the character which make it different from another. Feature extraction techniques differ from one application to another. Techniques that succeed in one application, may not be successful for other applications. Owing to the similar appearance of some different characters, the selection of the method for feature extraction remains the most important step in achieving high recognition accuracy. In this paper, DCT and DWT are used to assess their effectiveness in capturing the discriminative features of Arabic handwritten characters.

3.1. Discrete Cosine Transform

DCT is a technique to convert the pixel values of an image into its elementary frequency components [9]. By applying DCT on each character image, DCT coefficients of that image are obtained. It clusters high value coefficients in the upper left corner and low value coefficients in the bottom right of the array. Figure [2] illustrates the original image and that after DCT application. The DCT coefficients are extracted in a zigzag fashion and stored in a vector sequence. By applying DCT, each image of a character is represented by one vector. One of the characteristics of DCT is its ability to convert the energy of the image into a few coefficients [6]. By applying DCT on the character image with size 128×128 , 16,384 DCT coefficients of the image are obtained. The number of DCT coefficients chosen in the classification stage was set at 250 coefficients, rather than 16,384. Extensive experiments were carried out using different values of coefficients and it was found that 250 were the most appropriate. The number chosen was determined by empirical testing in order to reconstruct perceivable character. An overview of the extraction DCT coefficients is given in Figure [3].



Figure 2. Applying DCT on character image

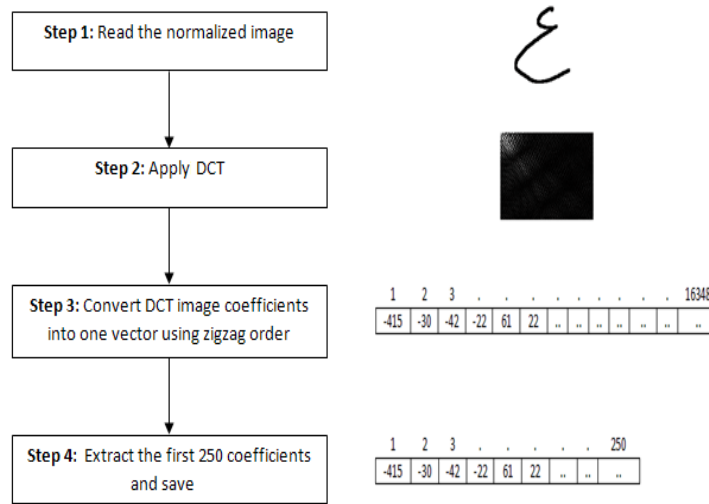


Figure 3. Overview of features extraction by DCT

3.2. Discrete Wavelet Transform

DWT is another technique used in this paper to extract the features of the characters where, at each decomposition level, a low-pass filter (LPF) and a high-pass filter (HPF) are applied to each row/column of the image to decompose this into one low-frequency sub-band (LL) and three high frequency sub-bands (LH, HL, HH) [11]. Figure [4] shows decomposition of DWT at one level. There are various types of wavelet transforms which can be applied, such as Haar and Biorthogonal, and others. Each of them has its particular features. From practical experiences, best results have been achieved in Arabic handwriting recognition by using Haar transform [12]. In this research, each character is decomposed by the Haar wavelet. The low frequency coefficients (LL) are close to the original image containing most details of the image and several works have used these coefficients to detect the features [11]. The image can be decomposed into more than one level, by repeating the processes on low pass filtering coefficients to provide the decomposition at level i . In this paper, each normalised character image is decomposed into three levels by the Haar wavelet and the low frequency coefficients (LL3) are used to extract the features of the characters. Therefore, by decomposition of a normalised image with a size of 128×128 into three levels, 256 low frequency coefficients are obtained. An overview of the extraction DWT coefficients at three levels is given in Figure [5].

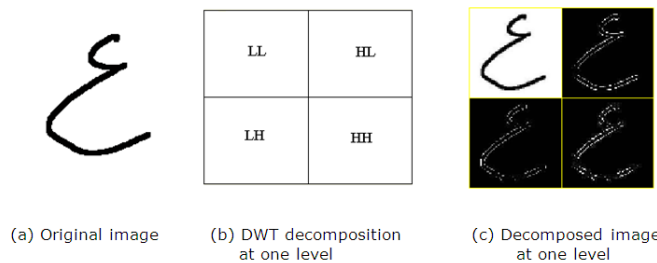


Figure 4. One level of decomposition of character image

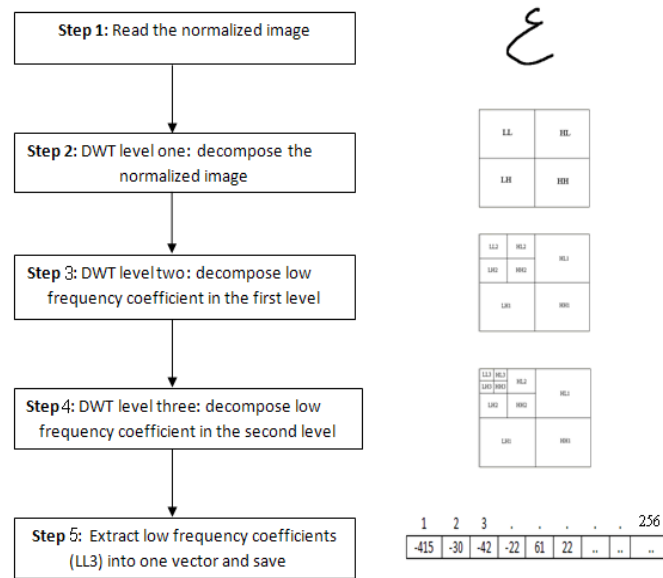


Figure 5. Overview of features extraction by DWT at three levels

4. Classification

A classifier is used to identify the shape of characters by using the features obtained by applying DCT and DWT, these are compared and saved as models for the trained classes. Features of an unknown shape of a character will be extracted and compared with the features of the training models to identify the unknown character shape. Common methods in the classification stage are the use of an ANN and HMMs [7]. Therefore, ANN and HMM are used to compare their effectiveness in classifying shapes of handwritten Arabic characters via the features captured by DCT and DWT.

4.1. Artificial Neural Networks

An ANN is a nonlinear system which is used widely in pattern classification [1]. It has been used to deal with the features that have been extracted from the shape of a character. An ANN consists of processing elements with weights which are learned from the training data. Three layers were used in this present research for the architecture of the network: the input layer, the hidden layer and the output layer. Figure 6 depicts an example of the architecture of the 3-layer ANN. The input layer is fed by the features of the shapes of characters. The number of nodes in this layer depends on the number of features extracted by DCT and DWT for each character. The last layer is called the output layer and the number of its nodes is based on the desired outputs. Therefore, the number of nodes in this layer depends on the number of shapes of characters. The hidden layer lies between the input and output layers. Feed-forward network multi-layer perceptron (MLP) back propagation (BP) with supervised training algorithm is used in this work. It is the best-known paradigm of training the ANN to classify patterns [13]. First, the ANN is trained by feeding it the features of the shapes of characters and the desired output. Then, it computes the errors between outputs and the desired outputs by using the values of the weights. These weights adjust to minimise the error between outputs and desired outputs. The testing phase is achieved by feeding

features of unknown characters to the network, which maps these features to the nearest character based on learnt characteristics. A classifier is used to identify the shapes of characters by using their features obtained by applying DCT and DWT. In experiments using DCT, 250 features are fed to the network as input signals, where in DWT 256 features are used. The number of nodes in the output layer depends on the number of character shapes. The number of nodes in the hidden layer was chosen experimentally to be 150 nodes to achieve the best performance.

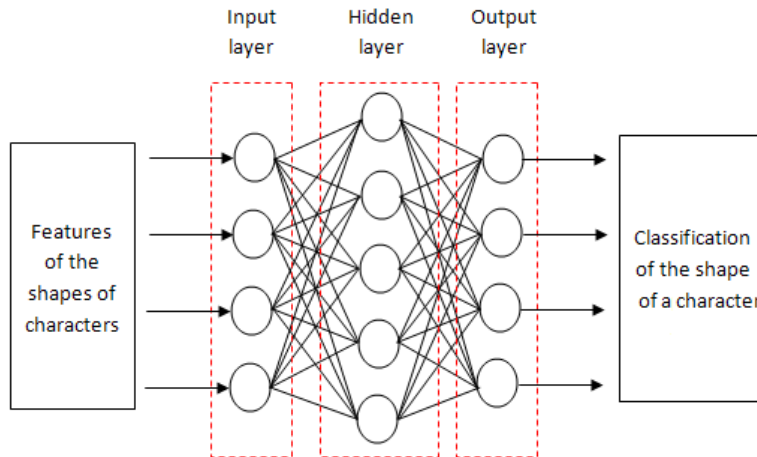


Figure 6. The architecture of the ANN used

4.2. Hidden Markov Model

An HMM is another technique used to deal with the features extracted from the characters. This technique is a very popular model in speech recognition. Owing to similarities between text and speech recognition, several researchers have used it for text recognition [4, 14]. HMM is used in this research to compare its effectiveness to that of ANN in classifying the shapes of handwritten Arabic characters via the features captured by DCT and DWT. There are many different topologies used in HMM based on classification. In this research, a left-to-right HMM is implemented. The sequence of state transition in both the training and testing of the model is related to the feature observations of Arabic handwritten characters. Figure 7 shows an example of 7 states for which a transition is allowed to the same state, the next state, and the following state only. HMM is used as another classification to identify the shape of characters by using their features obtained by applying DCT and DWT. The HMM toolkit is used to implement the HMM classifier. Each shape of character is represented once, by a 250-feature vector (DCT coefficients), and another by a 256-feature vector (DWT coefficients). States of various numbers were investigated and the best performance in representing the shape of character was selected. Although each character shape could have a different number of states, the same number of states was chosen for all the shapes.

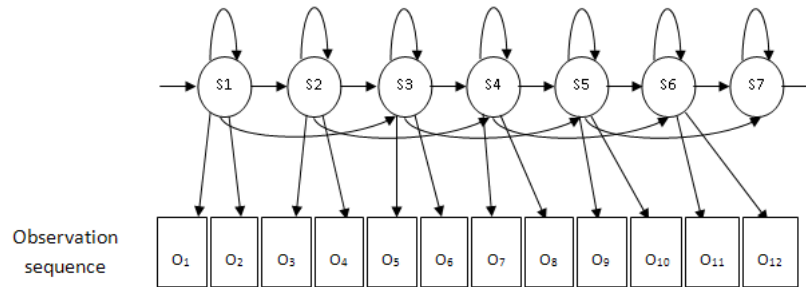


Figure 7. Example of a 7-state HMM

5. Comparative analysis of ANNs and HMM with DCT and DWT

Experiments were carried out using a dedicated HACDB database containing 6,600 shapes of Arabic handwritten characters which covered all the shapes of the Arabic characters, including overlapping characters considered as one shape. Two methods (DCT and DWT) were used and compared for the feature extraction process with ANN and HMM in the classification stage. To allow for a fair comparison, the size of images was set to 128×128 for both methods. The experiments were carried out in two steps. The first step was applied on the HACDB database containing 6,600 Arabic handwritten characters to classify them into 66 classes of individual shapes of Arabic handwritten characters, including overlapping characters. In the DCT technique, 250 coefficients were used to recognise the shape of a character. These coefficients were used by both ANN and HMM. On the other hand, in the DWT technique, 256 coefficients on level three were used, then used again by ANN and HMM. The second experiment was carried out by using the same database of Arabic handwritten characters which covered all the shapes of Arabic characters. Owing to the similarity between some of the handwritten characters, the number of classes was decreased to 34 by combining similar shapes of characters into the same class. For example, there are similarities between the shape of isolated characters and the shape of characters at the end of the word such as (ﺝ , ﺝ). Also, there are similarities between the shape of characters at the beginning of the word and the shape of characters in the middle of the word such as (ﺍ , ﺍ). Therefore, these similar shapes are grouped into the same class. The same number of DCT coefficients and DWT coefficients were used in the first and second experiment.

Table 1. Recognition rates by using different techniques

Methods	Recognition rate of the first experiment for 66 classes	Recognition rate of the first experiment for 34 classes
DCT & ANN	%75.31	%87.08
DWT & ANN	%59.85	%66.15
DCT & HMM	%73.92	%78.08
DWT & HMM	%54.08	%56.77

The results achieved in both experiments are summarised in Table 1. The second column in Table [1] shows the results of the first experiment obtained by using the different methods listed in the first column of the table for 66 classes, and the third column shows the results of the second experiment for 34 classes. Figure 8 illustrates

the comparison of performance of ANN with DCT and DWT, and HMM with DCT and DWT to recognise 34 classes. An increase in performance was noted in the second experiment when compared with the first one; this is due to some of the shapes being similar. Table [2] shows the performance of ANN with DCT when the original shapes of the characters are considered irrespective of the position of the characters. True positive rate and false positive rate for each class are also explained. The results also showed that the feature extraction by DCT with ANN yields a higher recognition rate. One reason is that most of the energy of the image is located in the upper-left corner of the 2D array, thus the DCT is more efficient, if the localisation of changes is significant [8, 15].

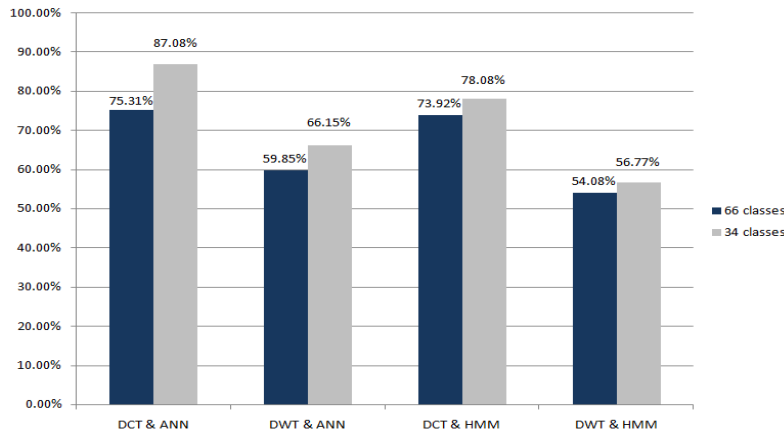


Figure 8. A comparison of ANN with DCT and DWT, and HMM with DCT and DWT

Table 2. Performance of ANN with DCT to recognise each class irrespective of the position

Classes	True positive rate	False positive rate
ا	95.00%	0.23%
ب	75.00%	1.84%
ح	77.25%	2.60%
د	75.25%	0.53%
ر	62.52%	1.07%
س	66.67%	3.84%
ك	72.50%	2.40%
ط	75.00%	1.70%
ع	77.50%	3.15%
ف	70.00%	1.76%
ق	62.50%	0.76%
ل	63.50%	0.85%
م	57.50%	0.92%
ه	81.00%	2.23%
و	75.00%	1.38%
ى	67.50%	1.38%
الح	70.00%	0.15%

لح	80.00%	0.23%
لمح	90.00%	0.15%
لا	93.33%	1.30%
لم	70.50%	0.15%
مح	80.25%	0.15%

Using only LL3 sub-band of DWT may affect the performance. However, the use of additional higher frequency sub-bands could improve the effectiveness of DWT. HMM is successful in classifying the words as it is based on the concept that it transforms the image into a sequence of observations. However, owing to the fact that the characters have fewer details than the words and the similarity between the characters, HMM seems less effective. The errors are mainly due to some individual handwriting styles, which make characters difficult to recognise. In some cases, the classification depends on the contextual information of the characters.

6. Conclusion and future work

This paper has compared the effectiveness of DCT and DWT in capturing the discriminative features of all shapes of Arabic handwritten characters with ANN and HMM in the classification stage. DCT and DWT techniques were used for feature extraction of the shapes of characters. Coefficients of both techniques were used with ANN and HMM for classification of the shapes of characters. The experiments were carried out in two steps. The first step was applied on the HACDB database containing 6,600 shapes of Arabic characters to classify them into 66 classes. The second experiment was carried out using the same database but to classify these characters into 34 classes. The results have demonstrated that the feature extraction by DCT with ANN yields a higher recognition rate. Although promising results have been achieved, there are some suggestions for further work to improve this performance:

- To increase the accuracy of recognition and reduce the number of similarities between the characters, Arabic characters can be divided into four groups: isolated characters; those at the beginning of the word; in the middle of the word; and at the end of the word. Each group is trained independently. This reduces the similarities between these groups and increases the rate of accuracy.
- Due to the size variations of characters and the similarity between some of handwritten characters, the structures of the features, such as loops curves and lines, can be extracted and classified into groups according to these features, this might improve the recognition performance.
- DWT decomposes an image into one low-frequency sub-band (LL) and three high frequency sub-bands (LH, HL, HH). In this paper, the low frequency coefficients (LL3) have been used to extract the features of the characters. Further work needed to investigation the performance when high frequency sub-bands are used. Other approaches, such as contourlet transform could used with DWT to improve the performance.

References

- [1] L. M. Lorigo and V. Govindaraju, "Offline Arabic handwriting recognition: a survey", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 28, no. 5, (2006), pp. 712–724.
- [2] D. Motawa, A. Amin and R. Sabourin, "Segmentation of Arabic cursive script", in *Proceedings of the 4th International Conference on Document Analysis and Recognition*, (1997), pp. 625–628, Washington, DC, USA.
- [3] S. Alma'adeed, C. Higgins and D. Elliman, "Recognition of off-line handwritten Arabic words using hidden markov model approach", in *16th International Conference on Pattern Recognition*, vol. 3, (2002), pp. 481–484.
- [4] R. El-Hajj, L. L. Sulem and C. Mokbel, "Arabic handwriting recognition using baseline dependant features and hidden markov modeling", in *Proceedings of the Eighth International Conference on Document Analysis and Recognition*, (2005), pp. 893–897, Washington, DC, USA.
- [5] S. Alma'adeed, "Recognition of off-line handwritten Arabic words using neural network", in *Geometric Modeling and Imaging-New Trends*, (2006), pp. 141–144.
- [6] J. H. AlKhateeb, R. Jinchang, J. Jianmin, S. S. Ipson and H. El-Abed, "Word-based handwritten Arabic scripts recognition using dct features and neural network classifier", in *5th International Multi-Conference on Systems, Signals and Devices*, (2008), pp. 1–5.
- [7] O. H. Assma, O. O. Khalifa and A. Hassan, "Handwritten Arabic word recognition: A review of common approaches", in *International Conference on Computer and Communication Engineering*, (2002), pp. 801–805.
- [8] A. Lawgali, A. Bouridane, M. Angelov, and Z. Ghassemlooy, "Handwritten Arabic character recognition: Which feature extraction method?", *International Journal of Advanced Science and Technology*, vol. 34, (2011), pp. 1–8.
- [9] A. Al-Haj, "Combined dwt-dct digital image watermarking", *Journal of Computer Science*, vol. 3, no. 9, (2007), pp. 740–746.
- [10] A. Lawgali, M. Angelova and A. Bouridane, "HACDB: Handwritten Arabic characters database for automatic character recognition", *4th European Workshop on Visual Information Processing (EUVIP)*, (2013), pp. 255 – 259.
- [11] M. Jiansheng, L. Sukang and T. Xiaomei, "A digital watermarking algorithm based on dct and dwt", In *Proceedings of the 2nd International Symposium on Web Information Systems and Applications*, (2009), pp. 104-107, Nanchang, China.
- [12] Z. Razak, N. A. Ghani, E. M. Tamil, M. Y. Idris, N. M. Noor, R. Salleh, M. Yaacob, M. Yakub and Z. B. Yuso, "Off-line jawi handwriting recognition using hamming classification", *Information Technology Journal*, vol. 8, no. 7, (2009), pp. 971-981.
- [13] A. K. Jain, M. Jianchang and K. M. Mohiuddin, "Artificial neural networks: a tutorial", *Computer*, vol. 29, no. 3, (1996), pp. 31–44.
- [14] S. Alma'adeed, C. Higgins and D. Elliman, "Off-line recognition of handwritten Arabic words using multiple hidden markov models", *Knowledge-Based Systems*, vol. 17, no. 2, (2004), pp. 75–79.
- [15] J. Dowling, B. Planitz, A. Maeder, J. Du, B. Pham, C. Boyd, S. Chen, A. Bradley and S. Crozier, "A comparison of dct and dwt block based watermarking on medical image quality", in *Proceedings of the 6th International Workshop on Digital Watermarking*, (2008), pp. 454-466.