# Handling Words Duplication and Memory Management for Digital Quran Based on Hexadecimal Representation and Sparse Matrix

**Ashraf Al-Omoush[1], Norita Md Norwawi[2,3]\*, Ahmad Akmalludin Mazlan[3]**

[1]*Faculty of Preparatory Year and Supporting Studies, Dammam, KSA*
[2]*Islamic Science Institute (ISI), Universiti Sains Islam Malaysia (USIM), Negeri Sembilan, Malaysia*
[3]*Faculty of Science and Technology, Universiti Sains Islam Malaysia (USIM), Negeri Sembilan, Malaysia*
*\*Corresponding author E-mail: norita@usim.edu.my*

## Abstract

Al-Quran is the holy book of the Muslims and the most important scripture containing knowledge on many domains. The recent advent of smart technologies like smart phones, digital devices and tablets has connected the daily life routines under a single touch adopted by many, these new tools with an exponential growth. This paper presented a Digital Quran Model (DQM) using hexadecimal representation using Unicode Hexadecimal and UTF-8 for character encoding, which is backward compatible with ASCII code. DQM target to handle all duplicated words or verses in Al-Quran using sparse matrix with double offset indexing to handle memory optimization. Three approaches were discussed: indexing and representation of the digital Quran to optimize storage, organize verses structure using sparse matrix to handle repetition with double offset indexing to efficiently use the space. The algorithms were implemented using Visual studio and Java server and the solution quality was measured by the size of a file before and after applying DQM model. For surah Al-Baqarah, the longest chapter in the Al-Quran, the reduction of the storage size was 25.00% whereas surah Al-Fatihah was 47.89%. The proposed DQM model is able to optimize the memory space and can be extended to other non-Roman characters used for information retrieval such as Hindi, Chinese and Japanese that are categorized in unicode standards.

*Keywords*: *Digital Quran; Hexadecimal Representation; Sparse Matrix; Unicode.*

## 1. Introduction

The popularity of digital Quran shown by the rapid increase of online Quran learners worldwide have resulted the needs of software applications that facilitate knowledge retrieval from the Quran, being the major source of authentic Islamic knowledge. In this study, we proposed Digital Quran Model (DQM) a representation based on hexadecimal with sparse matrix techniques to optimize memory storage. DQM handle all duplicated words or verses in the Quran by proposing a new technique for indexing based on hexadecimal representation that use Unicode Hexadecimal UTF-8 for character encoding which is backward compatible with ASCII code. In this study, surah Al-Baqarah and Al-Fatiha were taken as a test case to demonstrate how the technique is able to handle the many repeated or duplicated words in the Al-Quran.

This paper is organized as follows: Section 2 presents the review of related literature. Section 3 highlights the motivations and objectives of this study. The steps and methods of constructing the Digital Quran representation and algorithms based on hexadecimal and sparse matrix are presented in section 4, this is followed by section 5 presenting the implementation, and the evaluation is provided in section 6. Finally, our work in this paper is summarized in the last section.

## 2. Related Works

Al-Quran is the most important source of knowledge and considered as a standard text reference which relate valuable stories, prophetic habits and inherited Muslims wisdom. This great Book is a corpus which consists of 30 juzu, 60 hizb, 114 chapters, 6236 verses, 77439 words and 320015 letters [1]. Al-Quran is preserved from tampering since 14 centuries. This section presents the background and literature review of Al-Quran in the original Arabic text and digital Quran representation techniques that have been used.

Most recent notable work in the area of Quran studies using the cutting edge technology is by [2] and also their published work in [3, 4]. The main objective of this research deals with the design and development of a complete and comprehensive online cloud-based Quran portal. The portal and its applications makes all the reading and resource sections accessible to the audience whether the users are using laptops, PC, mobile, tablet, or personal digital assistants. Another notable work in the Quranic applications is presented in [5]. This work is about cross language information retrieval (CLIR). It presents a semantic technique on queries for retrieving more relevant results in CLIR, that concentrate on the Arabic, Malay or English query translation (a dictionary based method) to retrieve documents according to query translation. The size of the applications varies from one application to another with the highest size found to be 638MB and the smallest size of 79KB [6]. These applications are developed with the advancement of multimedia technology [7] where a primarily reciter of Al-Quran accomplish the learning of the holy book through a variety of sources, either websites such as Quran portals or modern gadgets like the iPhone and android apps as done in [8].

A survey related to the use of the internet or mobile devices to read the Quran or hadith were conducted in [1] where 73 % of the participants rely on the internet to find a particular Quranic verse or hadith and 52.9 % preferred reading the Quran in a soft copy or mobile devices as shown in Figure 1, whereas quite a large number of respondents use various types
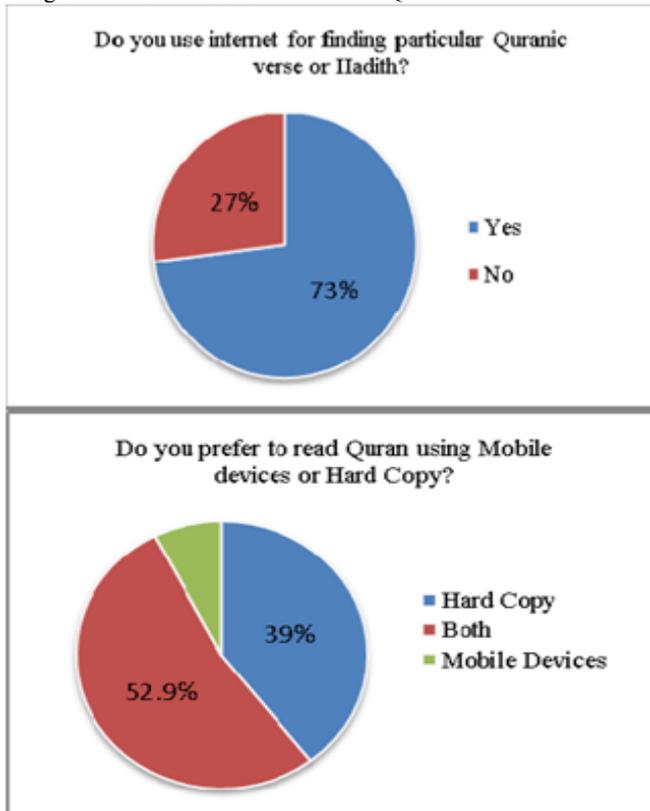
of digital devices to read and memorize the Quran.



**Fig. 1:** Use of Internet for reading Quran and preference of reading Quran using mobile devices [1]

Watermarking of text-documents has been classified as linguistic and non-linguistic by [9]. Linguistic techniques manipulate the lexical, syntactic and semantic properties of a document while trying to preserve the meanings; however, in the non-linguistic approaches changes are made to the text by using different text-attributes to embed a message. Text-watermarking techniques have been based on shifting techniques such as line-shift coding, word-shift coding and feature/character-coding, and natural-language based watermarking techniques such as synonym-substitutions or semantic-transformation techniques which are language-dependent. A study in [10] proposed a fragile watermarking method to preserve the authenticity of the digital Al-Quran. The method is considered as fragile watermarking method, which works on wavelet and spatial domains of digital Quran images. The authentication bits are embedded into each block of wavelet transformed image. Then, the least significant bits of pixels are considered to embed another authentication bits. The experimental results show that the watermarked image is imperceptible and fragile to the common attacks.

Works have also been done in the field of Digital Quranic Information Retrieval on different format and techniques but most researchers use the normal preprocessing techniques of Quran words and verses such as stemming, tokenizing, POS tagging and image processing. However, all of these techniques need time and storage, and ignore words duplication. Therefore, this paper presents a new technique which saves time and storage using UTF-8 characters encoding which is backward compatible with ASCII code [11], implemented using sparse matrix with double off set indexing to handle word duplication. Unicode transformation format (UTF) is the universal character code standard to represent characters where UTF-8 is an alternative coded representation form for all the characters in Unicode while maintaining compatibility with ASCII code [12].

## 3. Motivations and Objectives

The motivations behind this work come from the fact that all of our daily routine is being transferred to smart tools such as smart phone, tablet, PC and other devices in a similar way, people are transferring to smart devic-

es to read/brows their religious books. Millions of Muslims are also using smart devices to read/browse Quran and its sciences, which is easily available through the web or search engines. This inspired and motivated us to propose a new DQM that helps to optimized memory for the whole digital Quran by handling duplications. The purpose of this study is to propose DQM that can handle duplications and evaluate this implementation in terms of optimization of storage.

## 4. Methodology

The specific aim of this research work is to build a new model for digital Quran using hexadecimal representation and sparse matrix. The study begins with Arabic letters later develop the longest surah of the Al-Quran. By handling word duplication, it saves memory storage that may lead to ease of searching. Figure 2 shows DQM process flow. Letters, words, verses conversions and sparse matrix are discussed respectively in the next four sections as shown by the algorithm designed.
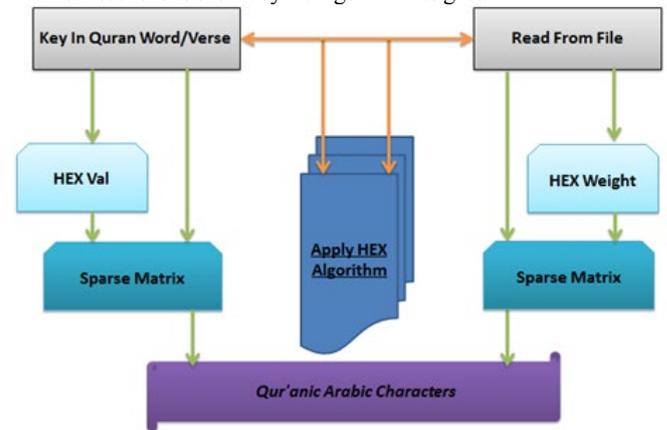


**Fig. 2:** DQM Process Flow.

### 4.1. Arabic Letters Conversion

Letters conversion is the backbone of the proposed technique using UTF-8 for character encoding, which is backward compatible with ASCII code. The placement of the letter in a word is significant which are: initial, medial, final and isolated such structure of Arabic letters make the process of digital Quran complex. The next example shows Thad (ض) letter in four expected positions and drawings on the word in addition to the hexadecimal representation that have been inherited from Unicode Standard 7.0, Copyright © 2014 Arabic Presentation, refer to Figure 3 and 4.

| isolated | ض | FEBD |
|----------|-----|------|
| initial | ضـ | FEBF |
| medial | ـضـ | FEC0 |
| final | ـض | FEBE |

**Fig. 3:** Thad (ض) letter in four expected position on the word.

| | | | | | | |
|---|---|---|---|---|---|---|
| C | FE7C | FE8C | FE9C | FEAC | FEBC | FEC0 | FED0 |
| D | FE7D | FE8D | FE9D | FEAD | FEBD | FEC1 | FED1 |
| E | FE7E | FE8E | FE9E | FEAE | FEBE | FEC2 | FED2 |
| F | FE7F | FE8F | FE9F | FEAF | FEBF | FEC3 | FED3 |

**Fig. 4:** Unicode Standard 7.0, Copyright © 2014 Arabic Presentation.

## 4.2. Words Conversion

The Al-Quran contains 77439 words and 18800 unique words, implying not only a huge amount of words but also repeated words. The new proposed technique represents the words of Al-Quran using Unicode by calculating the hexadecimal representation of each word. Transformed words are stored in an array. The word conversion is based on letter conversion as in the previous section constructed by the formula in Equation (1):

$$W(k) = L(h1, h2, \ldots hn) \qquad \text{(1)}$$
$$= \sum_{i}^{n} hi \qquad \text{(2)}$$

where *W(k)* is the hexadecimal value of the words and L is the function that calculates the value based on each letter representation in hexadecimal.

In addition, one must take into consideration the letter position on each word. The following example shows the weight of the word (قدير) QADEER in hexadecimal; according to Equation (1), calculated as follows:

W (قدير) = L (ق + ـد + ـي + ـر)
= L (FED7+FEAA+FEF3+FEAE)
= 3FB22

The size of the word (قدير) is 8 bytes compared to 5 bytes size of the same word in hexadecimal representation according to the above formula. This approach is useful in optimizing memory space, thus, may increase the speed in searching such as in information retrieval, text mining, Quran memorizing and interpretation.

## 4.3. Verses Conversion

Verses conversion depends on the previous section which is word conversion. The simple verse represented on hexadecimal is the first verse in the Quran بسم الله الرحمن الرحيم which is repeated 114 times; if we compare the size of the verse Bismillah Alrahman Alrahim ( بسم الله الرحمن الرحيم) will be 23 bytes compared to 4 bytes of FDFD hexadecimal representation, indicating storage size reduction.

Next, the case study of this research was the second chapter of the Al-Quran which was Surah Al-Baqarah ((البقرة)) which contains 26249 letters, 2279 unique words, 6140 words, 286 verses besides surah Al-Fatiha. Table 1 shows the Quran Surah Statistics (QSS) effectively made up 36 rows showing the Arabic printable letters, and 114 columns for all thesauruses. The rest of the rows or columns show computed cumulative statistics. The acronym WWR stands for Words without Repetition or unique words where the alphabets are ordered in accordance to the Unicode standard as shown in Table 1.

**Table 1:** Quran surat statistics (QSS) [1]

| Letter | Rank | Total | % | الفاتحة 1 | البقرة 2 | آل عمران 3 | النساء 4 |
|---|---|---|---|---|---|---|---|
| ء | 29 | 1578 | 0.48 | 0 | 125 | 68 | 71 |
| أ | 30 | 1511 | 0.46 | 0 | 126 | 74 | 64 |
| آ | 13 | 9119 | 2.76 | 1 | 657 | 359 | 433 |
| ؤ | 36 | 673 | 0.20 | 0 | 47 | 39 | 55 |
| إ | 18 | 5108 | 1.54 | 2 | 389 | 221 | 248 |
| ئ | 34 | 1182 | 0.36 | 0 | 94 | 64 | 54 |
| ا | 1 | 43542 | 13.17 | 23 | 3544 | 2005 | 2263 |
| ب | 9 | 11491 | 3.47 | 4 | 918 | 574 | 478 |
| ة | 24 | 2344 | 0.71 | 0 | 216 | 100 | 116 |
| ت | 10 | 10520 | 3.18 | 3 | 970 | 557 | 561 |
| ث | 31 | 1414 | 0.43 | 0 | 128 | 52 | 75 |
| ج | 21 | 3317 | 1.00 | 0 | 200 | 93 | 136 |
| ح | 20 | 4140 | 1.25 | 5 | 328 | 171 | 196 |
| خ | 23 | 2497 | 0.76 | 0 | 191 | 104 | 128 |
| د | 17 | 5991 | 1.81 | 4 | 458 | 252 | 301 |
| ذ | 19 | 4932 | 1.49 | 1 | 330 | 218 | 181 |
| ر | 8 | 12403 | 3.75 | 8 | 874 | 508 | 489 |
| ز | 28 | 1599 | 0.48 | 0 | 107 | 67 | 51 |
| س | 16 | 6012 | 1.82 | 3 | 451 | 227 | 306 |
| ش | 25 | 2124 | 0.64 | 0 | 168 | 86 | 82 |
| ص | 26 | 2072 | 0.63 | 2 | 155 | 88 | 122 |
| ض | 27 | 1686 | 0.51 | 2 | 133 | 66 | 101 |
| ط | 32 | 1273 | 0.38 | 2 | 99 | 50 | 65 |
| ظ | 35 | 853 | 0.26 | 0 | 62 | 36 | 45 |
| ع | 12 | 9405 | 2.84 | 6 | 797 | 383 | 404 |
| غ | 33 | 1221 | 0.37 | 2 | 75 | 62 | 62 |
| ف | 14 | 8747 | 2.64 | 0 | 751 | 396 | 503 |
| ق | 15 | 7034 | 2.13 | 1 | 553 | 306 | 255 |
| ك | 11 | 10497 | 3.17 | 3 | 832 | 485 | 582 |
| ل | 2 | 38191 | 11.55 | 22 | 3201 | 1892 | 1962 |
| م | 4 | 26735 | 8.08 | 15 | 2192 | 1246 | 1303 |
| ن | 3 | 27270 | 8.25 | 11 | 2019 | 1232 | 1334 |
| ه | 7 | 14850 | 4.49 | 5 | 1197 | 664 | 767 |
| و | 5 | 24813 | 7.50 | 4 | 2058 | 1147 | 1297 |
| ى | 22 | 2592 | 0.78 | 0 | 208 | 95 | 119 |
| ي | 6 | 21973 | 6.64 | 14 | 1596 | 998 | 1123 |
| Letters | | 330709 | | 143 | 26249 | 14985 | 16332 |
| WWR | | 14870 | | 26 | 2279 | 1469 | 1513 |
| Words | | 77797 | | 29 | 6140 | 3501 | 3763 |
| Verses | | 6236 | | 7 | 286 | 200 | 176 |

## 4.4. Sparse Matrix and Double off Set Indexing

Sparse matrix is the solution for words that has same hexadecimal representation by giving each word a unique ID to avoid duplication. This conversion is represented on a Lookup Table shown in Table 2. The conversion is represented with three columns; the first column contains the word in Arabic or Kalimah, the second column contains a counter that counts the number of each repeated word, the third column contains a hexadecimal representation for each word and the last column contains an ID. The storage of words being optimized through the use of one memory space for that particular word rather than one memory space for each Arabic character in the words of the Quran. Then, a sparse matrix will be constructed to encode to the display of Arab characters in Surah Al-Baqarah ((البقرة)) as shown in Figure 5. Here, the reduction of the storage was 25.00%.

**Table 2:** Lookup Table Represent Surat Al-Baqarah ((البقرة))

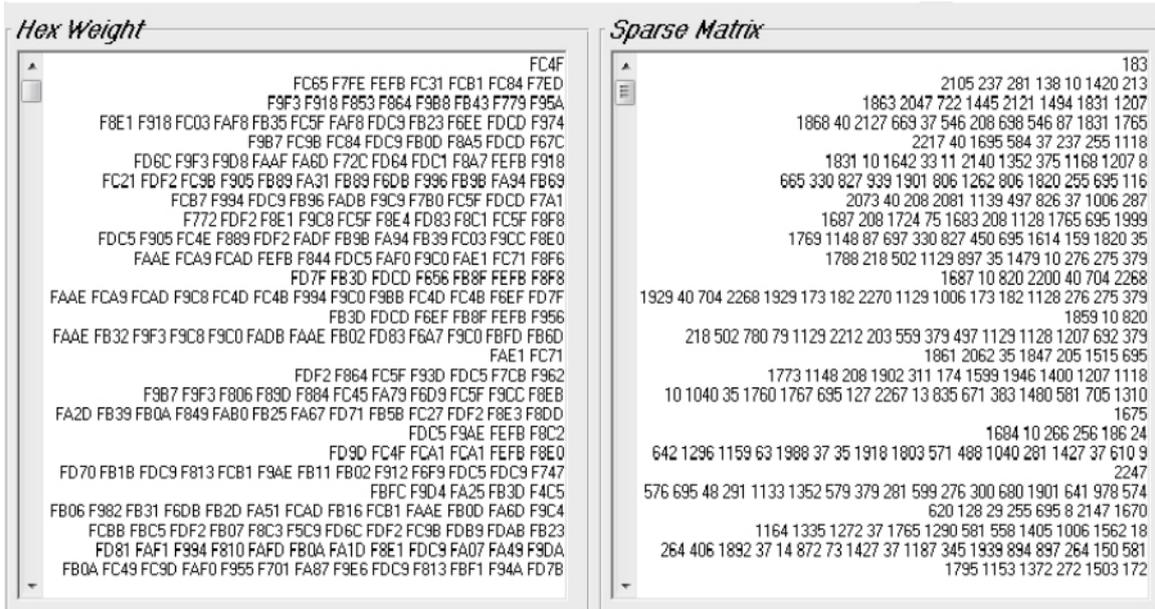| Kalimah | Count | New Hex | ID |
|---|---|---|---|
| إذ | 2 | FD32 | 1 |
| ذا | 2 | FD38 | 2 |
| حج | 2 | FD41 | 3 |
| رب | 2 | FD3C | 4 |
| تر | 2 | FD45 | 5 |
| شر | 2 | FD65 | 6 |
| آل | 2 | FD5E | 7 |
| إن | 2 | FD6C | 8 |

**Fig. 5:** Sparse Matrix Represent Arabic Characters for Surat Al-Baqarah

## 5. Implementation

This section discusses the implementation process of DQM from the first step until evaluation; developing a DQM prototype application involves combining and chaining many components, not all components will need to be executed on every system run. Instances containing the Quranic letters and words must be represented in Hexadecimal. Basically, the implementation process of the proposed system features three major processes which are: (1) letters conversion to hexadecimal, (2) word conversion to hexadecimal, and (3) sparse matrix with double offset indexing.

This Quranic Code will be done on three levels; character or letter, word, and verse level. Character level will be translated using UTF-8 character encoding for each of the Arabic characters in the Quran which is the basis of DQM, and compute the word representation in hexadecimal, followed by the verses. The algorithm was implemented using Visual studio and Java server as depicted in Figure 6 showing the main interface.
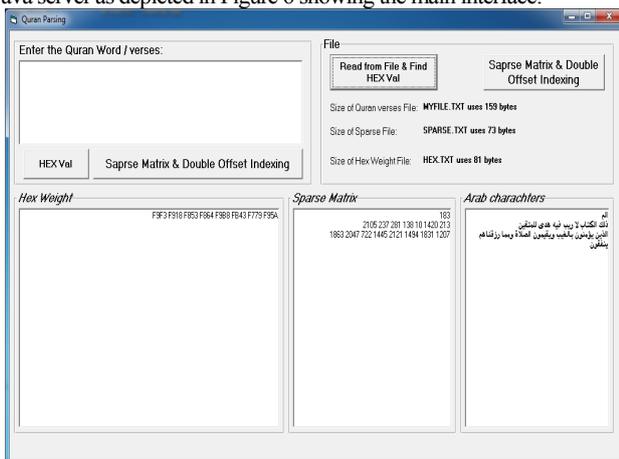


**Fig. 6:** Main interface of DQM.

In the DQM application, the algorithm can be applied directly by entering any Arabic or Quranic letter, word or verse, then retrieve the hexadecimal weight of the entered characters. Next, generate the sparse matrix and compare the size byte.

## 6. Evaluation

The main factor in the evaluation process is the solution quality, which is measured by comparing the size of a file before and after applying DQM algorithm. For example, the word Al-Rahim(الرحيم), Qadeer(قدير), and the first verse in the Quran which is *Bismillah Al-Rahman Al-Rahim* ( بسم الله الرحمن الرحيم), the hexadecimal representation was computed to be 5F893, 3FB22 and FDFD respectively which yield about 58.33%, 37.50%, 82.60% reduction.
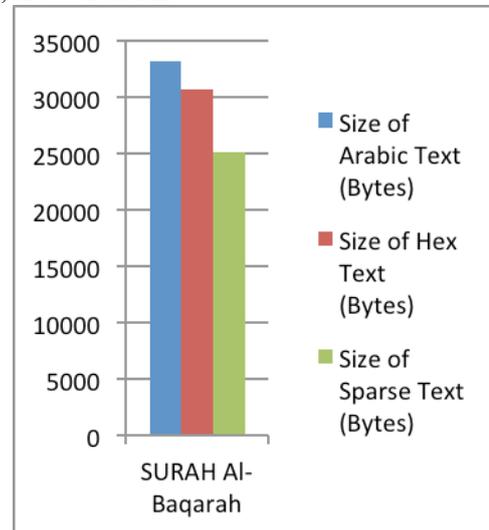


**Fig. 7:** Size comparison between Arabic text, Hex text and sparse text for Surah Al-Baqarah.

The first verse in the Al-Quran بسم الله الرحمن الرحيم is repeated 114 times. Therefore, if we compare the size of the verse Bismillah Al-Rahman Al-Rahim (بسم الله الرحمن الرحيم) will be 23 bytes instead of 4 bytes of FDFD hexadecimal representation. Imagine the reduction when repeated 114 times, the reduction will be of 2622 bytes compared to 456 bytes after applying the conversion algorithm which equal to 82.60%. For surah Al-Baqarah, the reduction on the storage size was 25.00%; and Al-Fatiha was 47.89%. Table 3 summarizes the comparison for all above mentioned words, verses, in addition to surah Al-Baqarah and Al-Fatiha. Figure 7 depicts the size comparison between arabic text, hex text and sparse text for surah Al-Baqarah.

**Table 3:** Summary of the comparison.

| Word/ Verse/ Surah | Number of Words | Repeated Words | Size of Arabic Text (Bytes) | Size of Hex Text (Bytes) | Size of Sparse Text (Bytes) | Occurrence in the Quran (Freq.) | Total Size Reduction (%) |
|---|---|---|---|---|---|---|---|
| Alrahim الرحيم | 1 | 0 | 12 | 5 | 9 | 148 | 58.33 |
| Qadeer قدير | 1 | 0 | 7 | 5 | 4 | 45 | 37.50 |
| Besm Ellah Alrahman Alrahim بسم الله الرحمن الرحيم | 4 | 0 | 23 | 4 | 4 | 114 | 82.60 |
| Surah Al-Fatiha الفاتحة | 29 | 3 | 159 | 149 | 85 | 1 | 47.89 |
| Surah Al-Baqarah البقرة | 2210 | 649 | 33199 | 30652 | 25149 | 1 | 25 |

# 7. Conclusion

In this paper, a new Digital Quran Model was proposed using Hexadecimal representation for storage optimization. Sparse matrix with double offset indexing use for handling words and verses duplications reduces the storage at 82.60% for the first verse in the Quran *Bismillah Al-Rahman Al-Rahim*, 25% for surah Al-Baqarah and 47.89% for Al-Fatiha. Thus, DQM is able to handle duplications and optimize the memory space. The approach can be extended to other non-Roman characters such as to serve information retrieval for non-English text (e.g. Hindi, Chinese, Japanese, etc.) categorized in unicode standards.

# Acknowledgement

# References

[1] Hakak, S., Kamsin, A., Tayan, O., Idris, M.Y.I, Gani, A. and Zerdoumi, S. (2017). Preserving Content Integrity of Digital Al-Quran: Survey and Open Challenges. *IEEE Access*, 5, 7305-7325.

[2] Drupal Open Source CMS. (n.d.). http://www.drupal.org.

[3] Adhoni, Z. A., Al Hamad, H., Siddiqi, A. A., & El Mortaji, L. (2013a). Towards a comprehensive online portal and mobile friendly Qur'an application. *Proceedings of the IEEE Taibah University International Conference on Advances in Information Technology for the Holy Quran and Its Sciences*, pp. 138-143.

[4] Adhoni, Z. A., Alhamad, H. A., Siddiqi, A. A.and Adhoni, E. Z. (2013b) CBQ-API: A Cloud-based Programming Interface for Quranic Applications. *Proceedings of the 3rd International Conference on IT Convergence and Security*, pp. 1-5, 2013.

[5] Yunus, M. A. M., Zainuddin, R., & Abdullah, N. (2013). Semantic method for query translation. *Int. Arab J. Inf. Technol.*, 10(3), 253-259.

[6] Khan, M. K., & Alginahi, Y. M. (2013). The holy Quran digitization: Challenges and concerns. *Life Science Journal*, 10(2), 156-164.

[7] Karkar, A., Alja'am, J. M., Eid, M., & Sleptchenko, A. (2015). E-Learning Mobile Application for Arabic Learners. *Journal of Educational and Instructional Studies in the World*, 5(2), 45-54.

[8] Adhoni, Z. A., & Siddiqi, A. A. (2013). A programming approach for the digital Quran applications. *International Journal of Engineering and Computer Science*, 13(5), 26-35.

[9] Adesina, A. O., Nyongesa, H. O., & Agbele, K. K. (2010). Digital watermarking: A state-of-the-art review. *Proceedings of the IEEE IST-Africa*, 2010, pp. 1-8.

[10] Kurniawan, F., Khalil, M. S., Khan, M. K., & Alginahi, Y. M. (2014). DWT+ LSB-based fragile watermarking method for digital Quran images. *Proceedings of the IEEE International Symposium on Biometrics and Security Technologies,* pp. 290-297.

[11] Al_Omoush, A., Norwawi, N. M., Ismail, R., Wahid, F. A., & Mazlan, A. A. (2017). Storage optimization for digital Quran using sparse matrix with hexadecimal representation. *Proceedings of the 6th International Conference of Computing and Informatics*, pp. 167-174.

[12] Diwakar, S., Goyal, P., & Gupta, R. (2010). Transliteration among Indian languages using WX notation. *Proceedings of the Conference on Natural Language Processing*, pp. 147-150.