

# Image compression using K-mean clustering algorithm

Amani Munshi , Asma Alshehri , Bayan Alharbi, Eman AlGhamdi, Esraa Banajjar, Meznah Albogami, Hanan S. Alshanbari

Department of Computer Science and Information System, Umm Al-Qura University, Makkah, Saudi Arabia  
s44286134@st.uqu.edu.sa, s44285259@st.uqu.edu.sa, s44285460@st.uqu.edu.sa, s44286654@st.uqu.edu.sa

## Abstract

With the development of communication networks, the processes of exchanging and transmitting information rapidly developed. As millions of images are sent via social media every day, also wireless sensor networks are now used in all applications to capture images such as those used in traffic lights, roads and malls. Therefore, there is a need to reduce the size of these images while maintaining an acceptable degree of quality. In this paper, we use Python software to apply K-mean Clustering algorithm to compress RGB images. The PSNR, MSE, and SSIM are utilized to measure the image quality after image compression. The results of compression reduced the image size to nearly half the size of the original images using  $k = 64$ . In the SSIM measure, the higher the K, the greater the similarity between the two images which is a good indicator to a significant reduction in image size. Our proposed compression technique powered by the K-Mean clustering algorithm is useful for compressing images and reducing the size of images.

**Keywords:** Types of images, RGB Images, Image Compression, K-mean Clustering, Peak Signal-to-Noise Ratio (PSNR), Mean Square Error (MSE), Structural Similarity Index Measure (SSIM).

## I. Introduction

Image compression plays a significant role in multimedia applications. Presently establishment of image compression is a type of data compression applied to digital images without degrading the quality of the image to an unacceptable level. The reduction in file size allows more images to be stored in a given amount of disk or memory space. It also reduces the time required for images to be sent over the Internet or downloaded from web pages.

We will be using the K-Means Clustering technique for image compression which is a type of transformation methods of compression. Using K-means clustering, we will perform quantization of colors present in the image which will further help in compressing the image.

The rest of this paper is organized as follows: - Section II, shows the related work. Section III, views the concept of images representation. Section IV describes the clustering technique. Section V gives the concepts of K-Means clustering .

In [10], the Singular Value Decomposition (SVD) technique for image compression and how to apply it is studied. This technique is based on dividing the image matrix into a number of linearly independent matrices. The

researcher used MATLAB to implement an image compression algorithm based on Singular Value Decomposition (SVD). Compression of medical images was the focus of the researchers' attention at [11]. Discrete Cosine Transform was used to compress images as was done by means of exploiting spectra similarity. Image quality evaluation is important in the case of image

Thereafter, section VI describes the methodology. Finally, section VII presents the results and discussion. Section VIII outlines the conclusion and future work directions. Documents are identified in italic within parentheses. Some components, such as multi-leveled equations, graphics, and tables are not prescribed, although the various table text styles are provided. The formatter will need to create these components, incorporating the applicable criteria that follow.

## II. Related work

In this part of the research, we will review previous literature working on image compression. Generally, images are compressed using multiple techniques. For example, Vector Quantization (VQ) and K-Means Clustering are commonly used to apply image compression. In research [6] a new algorithm (IDE-LBG) for generating optimum VQ Codebooks to efficiently compress grayscale images is proposed. This algorithm takes less computation time and results an excellent PSNR. The algorithm was tested on 5 different images, at a resolution of 512 x 512 and 6 different Codebook sizes. PSNR was calculated after each different compression. Based on VQ, the researchers presented in [7] a methodology for image compression. The methodology was tested on different image resolutions with a different Block Size. The MSE, PSNR (dB), and CR standards are used in this paper as performance measures. In [8] a scheme was presented to compress images with K-means clustering. The energy efficiency of the sensors was tested when sending data after applying a compression process. Energy consumption has been reduced by approximately 49% when sending images by the sensors. PSNR, MSE, SSIM was calculated as performance measures in this paper. Research [9] is concerned with medical images and trying to reduce the size of them to the lowest degree while preserving the quality of the images, as they are used in diagnosis. DWT-VQ (Discrete Wavelet Transform - Vector Quantization) technique is proposed for image compression. In the first

stage a preprocessing operation is used to remove the speckle and salt and pepper noises in ultrasound imaging. Then the proposed technique is applied to the image.

In [10], the Singular Value Decomposition (SVD) technique for image compression and how to apply it is studied. This technique is based on dividing the image matrix into a number of linearly independent matrices. The researcher used MATLAB to implement an image compression algorithm based on Singular Value Decomposition (SVD). Compression of medical images was the focus of the researchers' attention at [11]. Discrete Cosine Transform was used to compress images as was done by means of exploiting spectra similarity. Image quality evaluation is important in the case of image compression. There are several literatures that have established standards for image quality. For example, the SSIM scale was discussed in [12], which is used to measure the degree of similarity between two images. Also, there are two important metrics for evaluating image quality, MSE and PSNR [13].

### III. Images Representation

In computer science, images are represented in various forms, according to how the images are stored and the color data contained within them. Digital images can be encoded into two main approaches, Raster (bitmapped) images, and Vector images [1][2]. Fig. 1 explain the difference between raster and vector images. In vector graphics, points, lines, shapes, and polygons are used based on mathematical equations to represent this type of image. These images are mainly created by computer programs such as CAD and do not adopt resolution, so they can be reduced without losing the basic appearance. Most of the time these types of images are used in graphs. The vector graphics file formats can be SVG, EPS, PDF and AI [3].

In raster graphic images consist of a group of dots which are called pixel. These points are arranged in the form of a matrix with a number of columns and rows with a color for each point. Any type of image and color can be represented with this technique. The raster graphics file formats can be JPEG, PNG or GIF. There are several Color models that can represent colors [1]. Most of raster images can be saved according to two color models, CMYK and RGB [4]. In the RGB model, colors are displayed by combining the primary colors, red, green and blue, as shown in Fig. 2.

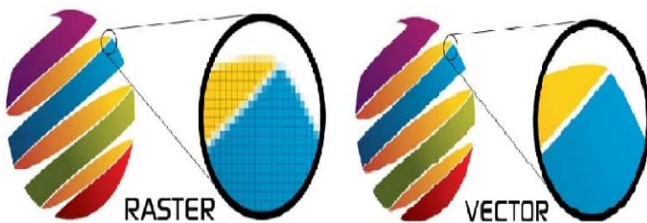


Fig. 1. Raster Images Vs Vector Images [3]

To represent each pixel color in the image three values for the red, blue, and green color must be set. As shown in Fig. 3, each basic color value can be represented by 8 bit = 1 bytes, so each pixel needs 3 bytes in the storage space to store it [5].

### IV. Clustering Techniques

Data clustering is a descriptive method for analyzing and grouping data [14]. There are many clustering algorithms that have been addressed by researchers in the literature. These algorithms can be classified into several categories. Common classifications of these methods are Partitioning methods, Hierarchical methods, Density-based methods, and Grid-based methods.

In partitioned clustering, a group of data objects is divided into non-overlapping groups. This technique is suitable for dividing the degree of colors in images, which enable easy compression. Also Partitioned clustering techniques can handle big data (images) more than other techniques. Under this classification there are several algorithms the most common of them is k-means. K-mean clustering is a simple unsupervised algorithm that can be applied to any form of data. also, it is easy to implement and its performance is very good since it is faster than other clustering algorithms [15].



Fig. 2. RGB Color Model [4]

NUMBERS			'RGB' = 3 SETS OF DIGITS		
R 255	R 102	R 51	11111111	01100110	00110011
G 0	G 102	G 204	00000000	01100110	11001100
B 0	B 255	B 153	00000000	11111111	10011001
R 255	R 255	R 51	11111111	11111111	00110011
G 255	G 0	G 204	11111111	00000000	11001100
B 102	B 204	B 255	01100110	11001100	11111111
R 51	R 51	R 255	00110011	00110011	11111111
G 51	G 51	G 153	00110011	00110011	10011001
B 0	B 153	B 153	00000000	10011001	10011001

Fig. 3. RGB image has three sets of numbers per pixel [5]

### V. K-means Clustering

Clustering is a process of dividing data into a group that shares certain characteristics according to patterns in the data. K-mean clustering is a clustering technique to identify clusters of data objects as explained in Fig. 4. Due to the nature of clustering, it is considered unsupervised learning, we do not need to put labels for the data as it recognizes some similar patterns between the data [16].

K-mean clustering is considered a centroid-based algorithm, whereby central points are chosen and the data close to each central point are grouped according to some properties in the data as explained in Fig. 5. A minimization of the sum of distances between the points and their respective cluster centroid is the main process of K-mean Clustering [17].

Distance metrics used in k-mean clustering algorithm to calculate the distance between data and centroids [18]. Euclidean Distance, Manhattan (City Block), Chebyshev Distance and Cosine Distance are the most common methods to compute the distance in the K-mean clustering algorithm. Euclidean Distance is measured using the following formula:

$$Dist_{xy} = \sqrt{\sum_{k=1}^m (x_{ik} - y_{ik})^2}$$

(1)

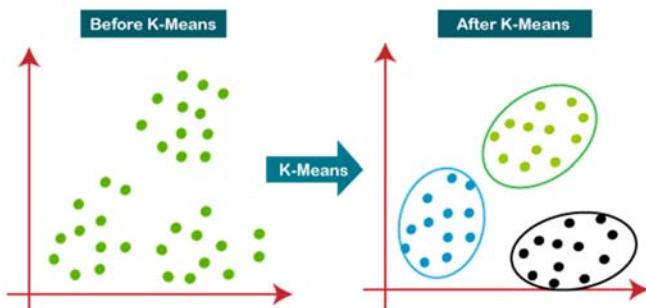


Fig. 4. K-Mean Clustering Technique [16]

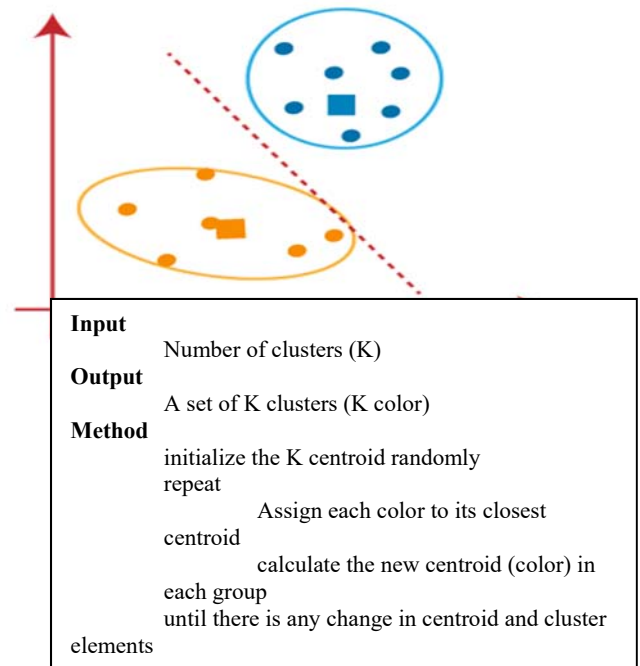


Fig. 5. K-Mean Clustering Technique [16]

### VI. Methodology

Digital images are 2- dimensional array of pixels, and they often require to be compressed in order to facilitate portability and storage [19]. K-means clustering is widely used to compress RGB images. Compression by K-mean clustering is a Lossy compression in which compressed images cannot be restored to their original state. However, the higher the compression ratio, the size decreases, but the compressed image quality will be affected. The pseudo code of the proposed methodology can be explained as follows:

#### A. Pre-processing

Before applying k-mean clustering, image data (pixels) must be prepared. Images are read from the storage media, then the data is read within (pixels with three bytes RGB color). These pixels are arranged in the form of a matrix with a number of columns and rows with a color for each point. The k-mean clustering needs to manipulate the image pixels as vector. The k-mean clustering needs to convert this array to a vector in order to facilitate handling of this data and facilitate setting of centroid points. The dimension of data is reduced from 2D to 1D. After K-mean clustering processing the dimension is again transformed into 2D.

#### B. Apply K-Mean Clustering Technique

As explained in section III, the image consists of pixels, each one of them is represented by 3 dimensions representing RGB intensity values, which ranges from 0 to 255. The

storage area for an image with dimensions of 600 \* 800 pixels can be calculated as follows:  $800 * 600 * 3 * 8 = 11520000$  bits. When using the K-mean clustering algorithm for compression of RGB images, a K-colors (centroid points) is selected. The rest of the colors are grouped into these centroid points according to the degree of similarity. The values of each pixel are replaced by the value of the centroid points. We can calculate the compression size for images as follows:

If k (centroid points) = 64 is placed, the image size will be  $600 * 800 * 6 + 64 * 3 * 8$ , where the third operand (i.e. number 6) represents the number of bits that can represent the values of k from 0 to 63. If the value of K is increased, the size of image will increase and the quality will increase. Finally, K-mean clustering algorithm is explained as follows:

- Image Input: Load the image from disk.
- Reshape Input Image: The input image must be changed from (rows, cols, 3), to (rows\*cols, 3).
- Clustering: Implement the k-Means clustering algorithm to find k-centroid points that represent its surrounding color combination.
- Replace each pixel with its centroid points: Replace the value of each of the pixels with its centroid point.
- Reshape Compressed Image: Again, reshape the compressed image to the original format (rows, cols, 3) dimensions.
- Output Compressed Image: Display the output image and store it to disk.

**C. Performance Metrics**

Image quality after compression is measured by using several standards metrics such as PSNR, MSE, and SSIM. Mean Square Error (MSE) gives the total amount of difference between two images. for  $I$  original image pixels, and  $I'$  is the compressed image pixels and the dimension for the two images  $M \times N$ , then the MSE can be computed as follows:

$$MSE = \frac{1}{M * N} \sum_{x=1}^M \sum_{y=1}^N (I(x, y) - I'(x, y))^2 \quad (2)$$

Peak Signal to Noise Ratio (PSNR) is a measure for image quality based on MSE, and can be computed as follows:

$$PSNR = 10 * \log \left( \frac{255^2}{MSE} \right) \quad (3)$$

Structural Similarity Index measure (SSIM) measures the similarity degree between two images. this measure can be computed based on luminance  $l$ , contrast  $c$  and structure  $s$  as follows:

$$SSIM(x, y) = [l(x, y)^\alpha * c(x, y)^\beta * s(x, y)^\gamma] \quad (4)$$

**VII. RESULT AND DISCUSSION**

In this paper Python software is used to apply k-mean algorithm to RGB images to compress them. The compressed images resulting from this process are nearly half the size of the original images using  $k = 64$ . Fig. 6 illustrate the Lena image before and after compression with  $k=32$ .



Fig. 6. Comparison Between Original and Compressed Image

Also, the compressed image quality is measured relative to the original. Based on MSE measure, the Peak Signal-to-Noise Ratio (PSNR) is commonly used to judge the quality of a compressed image relative to the original image. PSNR, MSE, SSIM performance measures are used in this research to compare the original and compressed image according to various K value. Table 1 and Fig (7-9) contain the comparison of the three metrics.

TABLE I. PERFORMANCE MEASURES COMPARISON FOR DIFFERENT K

Image	Lena			Baboon			Peppers		
	K=32	K=64	K=128	K=32	K=64	K=128	K=32	K=64	K=128
PSNR	32.15	34.3	36.2	27.17	29.19	31.13	29.34	31.68	33.71
MSE	39.66	24.16	15.61	124.9	78.29	50.11	75.66	44.2	27.66
SSIM	0.86	0.91	0.93	0.87	0.91	0.94	0.8	0.85	0.89



Compression Ratio	32%	48%	66%	40%	58%	74%	27%	41%	59%
-------------------	-----	-----	-----	-----	-----	-----	-----	-----	-----

especially with the development of the concept of IoT and the intense exchange of images between its elements.

### VIII.C ONCLUSION AND FUTURE WORK

In this paper a method for compressing images is proposed. The K-Mean clustering technique was used to implement this process. Previous literature concerned with

Fig. 7. Comparison MSE for different K and Images

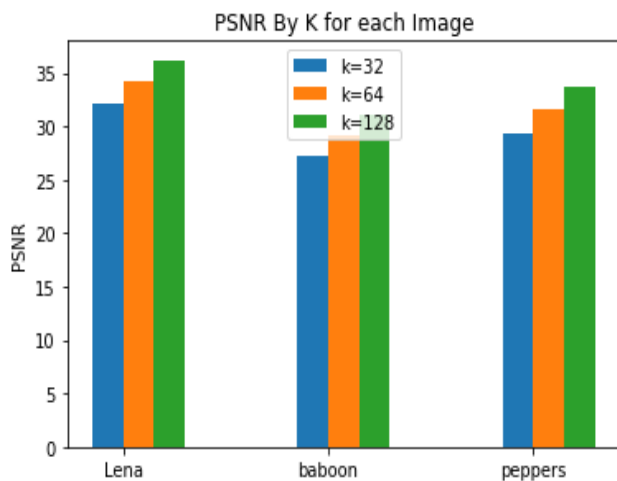


Fig. 8. Comparison PSNR for different K and Images

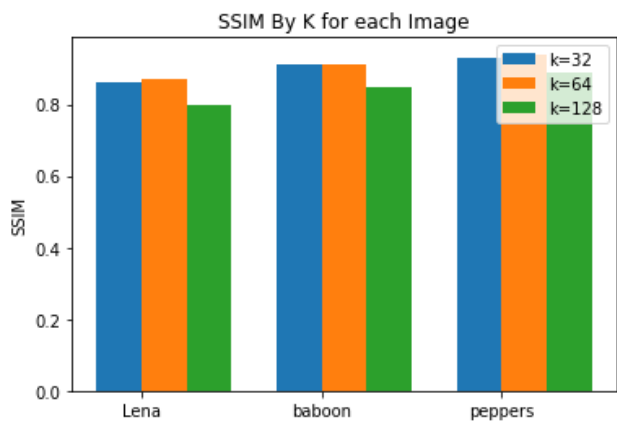


Fig. 9. Comparison SSIM for different K and Images

By increasing the quality of compressed image, the PSNR is increased, on the other hand the MSE decreased. SSIM is a measure of the extent of similarity between two images, so we see the higher the K, the greater the similarity between the two images, as shown in Table 1. For example, with a Compression Ratio of about 32%, the similarity is 86% between the original and compressed image, which is a very good ratio with a significant reduction in Image size. Finally, using the K-Mean clustering technique is very helpful,

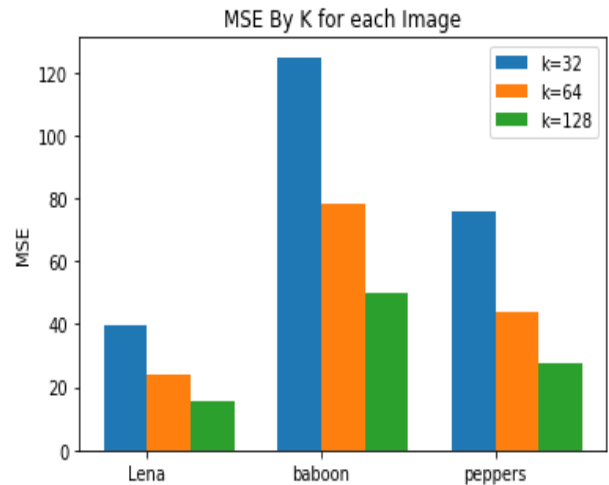


image compression has been explained. The K-Mean clustering technique is presented in this paper in general, then how to use it in image compression. The results demonstrate the benefits from compressing the images by reducing the size of the images while preserving the acceptable quality of the compressed images. For future work we intend to extend the images variety to include, measuring the quality of the proposed algorithm on high resolution images as well as, measuring the quality of night mode images. In addition we intend to expand the images format to include HEIF and HEVC.

### Acknowledgment

We would like to express our high regard to our families for their encouragement and inspiration supported us, and without which, we would not have come this far. Many thanks go to supervisor Dr. Hanan and our deep appreciation for continuous guidance and her prompt help and provide advice support to helped us finalize our project and offered deep insight into the study. Also, special thanks should be given to group friends that worked on this project for the kindness, cooperation, positive energy, constant motivational words and caring throughout the whole project.

### References

[1] Gouet-Brunet V. (2009) Image Representation. In: LIU L., ÖZSU M.T. (eds) Encyclopedia of Database Systems. Springer, Boston, MA. [https://doi.org/10.1007/978-0-387-39940-9\\_1438](https://doi.org/10.1007/978-0-387-39940-9_1438)

- [2] Marques O. (2008) Image Data Representations. In: Furht B. (eds) Encyclopedia of Multimedia. Springer, Boston, MA. [https://doi.org/10.1007/978-0-387-78414-4\\_343](https://doi.org/10.1007/978-0-387-78414-4_343)
- [3] Claudia Jeffrey, Raster vs Vector Graphics – Ultimate Guide, May 21, 2020, accessed March 2021.
- [4] Applying Color Theory to Digital Media and Visualization, Rhyne, Theresa-Marie, Boca Raton, FL: CRC press, 2016. 184 pp. ISBN 9781498765497.
- [5] Image Pixels, <http://shutha.org/node/789> accessed in March 2021.
- [6] Nag, S. (2017). Vector Quantization using the Improved Differential Evolution Algorithm for Image Compression. ArXiv, abs/1710.05311.
- [7] Adokar, D. U., & Gurjar, A. A. (2020). Image Compression using Vector Quantization. Grenze International Journal of Engineering & Technology (GIJET), 6(2), p. 69.
- [8] Paek, J., & Ko, J. (2017). K-Means Clustering-Based Data Compression Scheme for Wireless Imaging Sensor Networks. IEEE Systems Journal, 11, 2652-2662.
- [9] Ammah, P.N., & Owusu, E. (2019). Robust medical image compression based on wavelet transform and vector quantization. Informatics in Medicine Unlocked, 15, 100183.
- [10] K. Mounika, D. Sri Navya Lakshmi, K. Alekya and M.R.N. Tagore, “SVD Based Image Compression”, International Journal of Engineering Research and General Science, Vol. 3, No. 2, pp. 1-5, 2015.
- [11] Wu, Y.G. and S.C. Tai, 2001. Medical image compression by discrete cosine transform spectral similarity strategy. IEEE T. Informat. Technol. Biomed., 5(3): 236-243.
- [12] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, “Image quality assessment: from error visibility to structural similarity,” Image Processing, IEEE Transactions on, vol. 13, no. 4, pp. 600–612, 2004
- [13] Z. Wang and A. Bovik, “Mean squared error: Love it or leave it? a new look at signal fidelity measures,” Signal Processing Magazine, IEEE, vol. 26, no. 1, pp. 98–117, Jan. 2009
- [14] J. HAN AND M. KAMBER, Data mining: concepts and techniques, Morgan Kaufmann Publishers, Inc., 3<sup>rd</sup> edition 2011.
- [15] D. Lam and D. C. Wunsch, “Clustering,” Academic Press Library in Signal Processing,” Signal Processing Theory and Machine Learning, vol. 1, 2014.
- [16] K-Means Clustering Algorithm, <https://www.javatpoint.com/k-means-clustering-algorithm-in-machine-learning>, Accessed in March 2021.
- [17] D. Sisodia, L. Singh, S. Sisodia, K. Saxena, “Clustering Techniques: A Brief Survey of Different Clustering Algorithms”, International Journal of Latest Trends in Engineering and Technology (IJLTET), Vol.1 Issue3 September 2012.
- [18] Arzoo, K., & Rathod, K.R. (2017). K-Means algorithm with different distance metrics in spatial data mining with uses of NetBeans IDE 8.2.
- [19] Mr. Chandresh K Parmar, Prof. Kruti Pancholi,—A Review on Image Compression Techniques! Journal of Information, Knowledge And Research in Electrical Engineering ISSN:0975–6736 volume–02, Issue–02 Nov12 to Oct13