

# Analysis of Image Retrieval Accuracy for Searching Medicinal Plant Information

Jumi

Department of Bussiness Administration  
Politeknik Negeri Semarang,  
Semarang, Indonesia

Azizah

Department of Bussiness administration  
Politeknik Negeri Semarang,  
Semarang, Indonesia

Tedjo Mulyono

Department of Civil Engineering  
Politeknik Negeri Semarang,  
Semarang, Indonesia

Nur Hayati

Department of Information Technology  
Institut Teknologi dan Sains Nahdlatul Ulama Pekalongan,  
Pekalongan, Indonesia

**Abstract**—Image retrieval system is an alternative to image-based information search keys. Information searches do not only use text as a search key but can also use images. The use of images as search keys must have similarities or similarities. Image similarity affects the accuracy of the search result information. In this research, the accuracy analysis of image similarity is carried out using texture, color and shape features in medicinal plant leaf images. Image similarity accuracy analysis is measured using the similarity of texture, shape, color features and the combination of the three features. The feature extraction method uses statistical texture for texture features, color moment for color features and shape features using slimmness, roundness, rectangularity, and narrow factor. Before feature extraction, the image preprocessing stage is first carried out. This stage uses the Grayscale, resize, Histogram Equalization and edge Enhancement methods. Then the results of feature extraction are stored in a database that has gone through the clustering stage using K-Means. The results showed that the accuracy of image retrieval with test data of 300 images of medicinal plant leaves, reached more than 90% in the weighting scheme  $W_s$  (weighted Shape) = 40% and  $W_c$  (weighted color) = 40%,  $W_t$  (Weighted Texture) = 20%. In Clustering 15, the average time for the image retrieval process is less than 5 milli-seconds.

**Keywords:** leaf, image, feature, retrieval, statistical.

## 1. INTRODUCTION

Medicinal plants thrive with various types and various properties in Indonesia. The variety of medicinal plants is very diverse and has different leaf characteristics including texture, shape and color. The wealth of medicinal plants in Indonesia has the potential to be managed because it has economic value that supports the income of the community or country. There are more than 30,000 types of plants in Indonesia and more than 2000 are herbal medicinal plants [1]. The management of medicinal plants with a large number of species requires a system to make it easier for the public or drug industry players

to get information about these medicinal plant species. At this time people still rely on chemical drugs because of the lack of information related to the benefits or properties of herbal medicinal plants. The availability of a medicinal plant information system that is easily accessible and has the right accuracy about the benefits of medicinal plants will reduce people's dependence on chemical drugs. Information obtained can be in the form of digital images which are then analyzed and processed by the system. Information on medicinal plants can be done by recognizing the leaves. The system searches for herbal medicinal plants and recognizes a pattern or characteristic of the object and is then used for the retrieval key to search for medicinal plant information. Research on image retrieval systems has long been developed, one of which is by distinguishing the texture of the image. Image texture can be distinguished by density, regularity, uniformity, and roughness [2]. Because computers cannot distinguish textures like human vision, texture analysis is used to determine the pattern of a digital image. Texture analysis will produce a value of texture characteristics that can then be processed by a computer for the classification process [3].

Medicinal plants are one type of plant that has a very important role in human life, in addition to the stem, flowers, leaves, fruit, and roots that are most often used in the identification process of medicinal plants, namely leaves. However, physical characteristics in the form of color are considered not so significant in determining the type of leaf. This is because almost all types of leaves have a dominant green color. Meanwhile, to obtain shape features, there are sometimes difficulties in capturing whole leaf data, especially for leaves that have a large scale, so that the combination of shape, color and texture features of the leaves is not significant.

is a more appropriate feature used in the identification of herbal medicinal plant leaves. The leaves that will be studied in this study are the types of leaves of traditional medicinal plants that include herbal plants. In this study, researchers want to create a system that can identify traditional medicinal

plants found around to help lay people in the identification process. The identification process is carried out by recognizing the shape, color and texture of the leaf image worn on a type of traditional medicinal plant leaf and matching it with data from the medicinal plant leaf image database. The recognition process is carried out by extracting features using the Invariant moment, Color Moment and Statistical Texture methods. The features of each leaf are used for the identification process using the Euclidean Distance and K-Means Clustering methods. This research aims to recognize the type of leaves of medicinal plants based on shape, color and texture in determining the type and efficacy of medicinal plants.

This research aims to analyze image-based rice type identification using feature value weighting scheme and computation time calculation after clustering using K-Means technique. The K-Means algorithm is used to determine the cluster position on each image by first calculating the image distance with all centroids using the euclidean distance method.

## 2. RELATED WORKS

### 2.1 Statistical Texture

Region description is an important approach to calculate its texture content. There are three principle approaches to describe the texture of a region, namely statistical, structural and spectral [4]. Statistical is one of the approaches to describe texture by using statistical moments from the intensity histogram of an image or region.

#### a) Smoothness with Mean Gray level

If a random variable  $z$  denotes intensity, and  $p(z_i)$  is its histogram with  $i = 0, 1, 2, \dots, L-1, L$  being the number of intensity levels, then the  $n$ th moment of  $z$  with respect to the mean is in equation (1).

$$\mu_n(z) = \sum_{i=0}^{L-1} (Z_i - m)^n p(Z_i) \tag{1}$$

With  $m$  being the mean  $Z$  value (mean intensity), which is shown by equation (2).

$$m = \sum_{i=0}^{L-1} Z_i p(Z_i) \tag{2}$$

Based on equation (1) and equation (2), the values  $\mu_0=1$  and  $\mu_1=0$  are the values of the 0th and 1st moments. The second moment or  $\mu_2$  called the variance  $\sigma^2(Z)$  is an important part of texture description. It is a measurement of intensity contrast that can provide a description of *relative smoothness*. The variance value is calculated using equation (3) [5]

$$\sigma^2 = \sum_{i=0}^{L-1} (Z_i - m)^2 P(Z_i) \tag{3}$$

#### a. Deviation Standard

Standard deviation, denoted by  $\sigma(z)$ , is also often used as a texture measurement.

$$\sigma_j = \sqrt{\left( \frac{1}{L} \sum_{i=0}^L (P_{ji} - m_j)^2 \right)} \tag{4}$$

#### b. Skewness

$$\mu_3(Z) = \sum_{i=0}^{L-1} (Z_i - m)^3 p(Z_i) \tag{5}$$

Skewness is a measurement of the skewness of the histogram while the fourth moment  $\mu_4$  is a measurement of relative flatness.

#### c. Uniformity

Another texture measurement that can be used as a description of an image based on its histogram is the uniformity measurement, there are:

$$U(z) = \sum_{i=0}^{L-1} P^2(Z_i) \tag{6}$$

#### d. Entropy

As for measuring an average entropy, use Equation (7).

$$e = - \sum_{i=0}^{L-1} P(Z_i) \log_2 P(Z_i) \tag{7}$$

The Variable  $p$  has a value range of  $[0,1]$  and the sum is equal to one, the  $U$  measurement is the maximum for an image with all gray levels the same (maximum uniformity), and decreases from this maximum value. Entropy is a measurement of change (variability) and is zero for a constant image.

Statistical Moment Distance Value To calculate the statistical moment distance value, use equation (8)

$$D_{\text{Tekstur}}(H, I) = \sum_{i=1}^r W_{i1} |m_i^1 - m_i^2| + W_{i2} |\sigma_i^1 - \sigma_i^2| + W_{i3} |\mu_i^1 - \mu_i^2| + W_{i4} |U_i^1 - U_i^2| + W_{i5} |e_i^1 - e_i^2| \tag{8}$$

### 2.2 Color Moments

In general, there are three basic colors, where each color can be reproduced by combining or mixing a set that corresponds to the three basic colors. The number of colors is influenced by three color component vectors in three-dimensional space in a coordinate system. The set of all colors forming a vector space is called a color space or color model.

The results of light perception in the spectrum of the region visible to the retina of the eye are colors with wavelengths between 400nm and 700nm. RGB colors have the visualization of a cube that has 3 axes to represent the colors R (red) or red, G (green) or green and B (blue) or blue [6]. One of the opposite base corners represents black when the value  $R=G=B=0$ , while the opposite top corner represents white when the value  $R=G=B=255$  (8 bit color system). Color can be described as a metric representation of color space so that color differences can be calculated using distance differences via Euclidean distance. The transformation equation shown by equations 9 to 11 is used to convert vector values from RGB to HSV.

$$r = \frac{R}{(R+G+B)}, g = \frac{G}{(R+G+B)}, b = \frac{B}{(R+G+B)} \tag{9}$$

$$V = \max(r, g, b)$$

$$S = \begin{cases} 0, & \text{jika } V = 0 \\ 1 - \frac{\min(r,g,b)}{V}, & V > 0 \end{cases} \tag{10}$$

$$H = \begin{cases} 0, & \text{jika } S = 0 \\ \frac{60*(g-b)}{S*V}, & \text{jika } V = r \\ 60 * \left[ 2 + \frac{b-r}{S*V} \right], & \text{jika } V = g \\ 60 * \left[ 4 + \frac{r-g}{S*V} \right], & \text{jika } V = b \end{cases} \tag{11}$$

$$H = H + 360 \text{ jika } H < 0$$

Color uses three main moments from the image color distribution, namely mean, standard deviation, and skewness, so that this method produces three values for each color component [7]. Each color component, namely HSV (Hue, Saturation and Value) has 3 moments. Calculation of these three moments uses equation 12, equation 13 and equation 14.

The results of preprocessing using histogram equalization shown in Figure 2 show that the image has more contrast. This method is able to differentiate images based on color features [8].

Color Moment is a fairly good method for recognizing color characteristics. This method uses three main moments of image color distribution, namely mean, standard deviation, and skewness [9]. Each color component, namely HSV (Hue, Saturation and Value) has 3 moments. Calculation of these three moments uses equation (12), equation (13) and equation (14).

a. Moment 1 – Mean :

$$E_i = \sum_{N=1}^{j=1} \frac{1}{N} P_{ij} \tag{12}$$

Mean: can be said to be the average color value in the image.

b. Moment 2 – Varians :

$$\sigma_i = \sqrt{\frac{1}{N} \sum_{N=1}^{j=1} (P_{ij} - E_i)^2} \tag{13}$$

Standard deviation: the range of spread of data from the mean

c. Moment 3 – Skewness :

The measure of data asymmetry around the mean is called skewness which is calculated by equation (14).

$$S_i = \sqrt[3]{\frac{1}{N} \sum_{N=1}^{j=1} (P_{ij} - E_i)^3} \tag{14}$$

In Equation (15), the distance value from the color distribution of the query image can be calculated

$$D_{color} (H, I) = \sum_{i=1}^r W_{i1} |E_i^1 - E_i^2| + W_{i2} |\sigma_i^1 - \sigma_i^2| + W_{i3} |S_i^1 - S_i^2| \tag{15}$$

E = Average image color value (Mean).

- (H,I) = Two images are compared
- i = HSV color Component Index(H=1, S=2, V=3)
- r = Number of Indexes ( 3 )
- σ = The square root of the variance (Standard Deviation).
- S = A measure of the degree of asymmetry (Skewness).
- N = The total number of pixels in the image.
- J = Pixel order.
- Wi = The weight of each moment.
- Pij = The value of the i color component at the j pixel

### 2.3 Shape of Feature

There are four types of features used to determine the similarity of leaf shapes. There are 4 features used to determine the similarity of leaf image shapes there are [10]:

a) Slimness

Slimness is the ratio between leaf length and leaf width [11].

$$Slimless = \frac{Lp}{Wp} \tag{16}$$

where Lp is the length of the leaf, and Wp is the width of the leaf.

b) Form Factor / Roundness

$$Roundness = \frac{4xA}{p^2} \tag{17}$$

where A is the area/area of the leaf and P is the circumference of the leaf.

c) Rectangularity

Rectangularity R describes how similar a leaf shape is to a rectangle based on the ratio of the area A to the minimum rectangular boundary LW [12]. Rectangularity calculation is carried out using Equation (17).

$$R = \frac{A}{LW} \tag{17}$$

d) Narrow factor

Narrow factor NF is the ratio of the leaf diameter D to the length of the major axis L [13]. Narrow factor is calculated with Equation (18).

$$NF = \frac{D}{L} \tag{18}$$

### 3. METHOD

Image retrieval testing in this study uses a method by combining texture, color and shape features. Testing uses three stages, namely data acquisition through object capture, image preprocessing, image feature extraction, database recording, clustering, matching similarity between input images and images stored in the database for recognition and identification shown in Figure 1. This research uses a combination method of texture, shape and color feature values then weighting the three feature values and clustering that is not found in previous research, namely a recognition method based on texture, color and shape features for the retrieval process in the medicinal plant database using leaf image search keys.

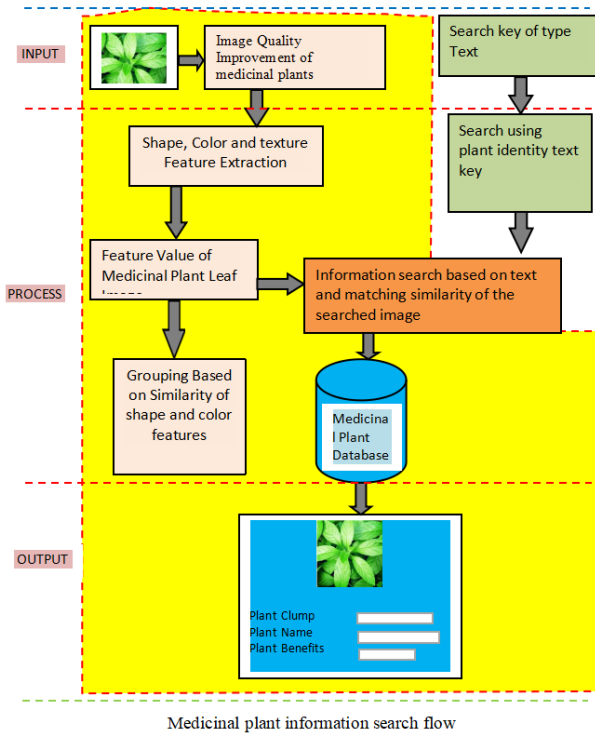


Fig. 1. Overview of leaf image retrieval

The details of the stages of this research are as follows:

### 3.1 Preprocessing

Preprocessing is a stage of improving image quality before feature extraction with the aim of increasing the accuracy of image feature extraction results. There are differences in preprocessing in color feature extraction and shape feature extraction. The difference in preprocessing is to get a quality image before feature extraction.

#### Color Feature Preprocessing

##### a) Resize

The larger the image size, the longer the extraction time, so a resize stage is needed to speed up the computation process. At this stage, the image is resized to 200 x 200 pixels.

##### b) Histogram Equalization

At this stage, histogram smoothing is carried out so that the image quality becomes more contrasting to get quality color feature values. The results of histogram smoothing can be seen in Figure 2.

The preprocessing process of medicinal plant leaf images carried out in this study can be seen in figure 2.

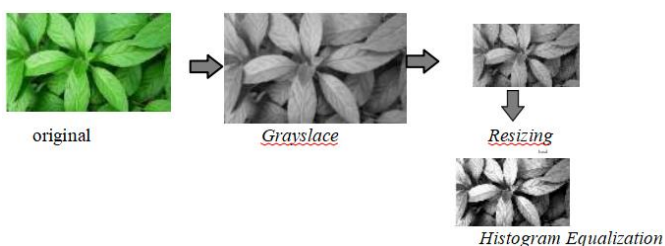


Fig. 2. Preprocessing Stages of Histogram

### 3.2 Feature extraction

Extraction of texture, color and shape features is performed on the image after going through the preprocessing stage. Texture feature extraction using statistical texture, color feature using color moment and shape feature extraction using slimmness, roundness, rectangularity, narrow factor.

In statistical texture smoothness, Standard Deviation, uniformity, Skewness and entropy. While Color moments use three main moments there are mean, standard deviation, and skewness, so this method produces three values for each color component [14]. The Color Moments method is a method used to distinguish images based on their color features [15]. Color moments assume that the color distribution in an image can be expressed as a probability distribution. Extraction of shape features using slimmness, roundness, rectangularity, narrow factor

### 3.3 Feature Weighting

Image retrieval accuracy has a different percentage for each combination of features used for testing. The use of various feature value weighting models for retrieval accuracy testing is done to get a combination of features that have an accuracy of more than 85%. Weighting is done with four variations, namely the W1 scheme (0.40, 0.30, 0.30) is a 40% weighting scheme for the shape feature weight, 30% color feature and 30% texture feature weight. Scheme W2 (0.6, 0.20, 0.20) is a weighting scheme of 60% shape feature weight, 20% color feature and 20% texture feature weight. Scheme W3 (0.50, 0.25, 0.25) is 50% weight of shape features, 25% color features and 25% weight of texture features. And the W4 scheme (0.40, 0.40, 0.20) is a 40% weighting scheme for shape feature weights, 40% color features and 20% texture feature weights.

### 3.4 Clustering

Data clustering to speed up the identification process. In this research, the K-Means Clustering method and Euclidean distance are used as image similarity measurements[16]. The similarity measurement equation uses equation 20.

$$D(Q, M) = \sqrt{\sum_{n=1}^k (Q_n - M_n)^2} \quad (20)$$

Where  $Q$  and  $M$  are the features of the query image and the database image in the  $n$ th dimension.

Image similarity is determined from the difference between the feature values of the identified image and the image in the database. Feature value differences that are close to zero have the highest level of similarity

4. RESULT AND DISCUSSION

Testing in this study used 10 species of medicinal plants with a total of 300 rice images measuring 200 x 200 pixels. Sample testing data can be seen in Figure 3

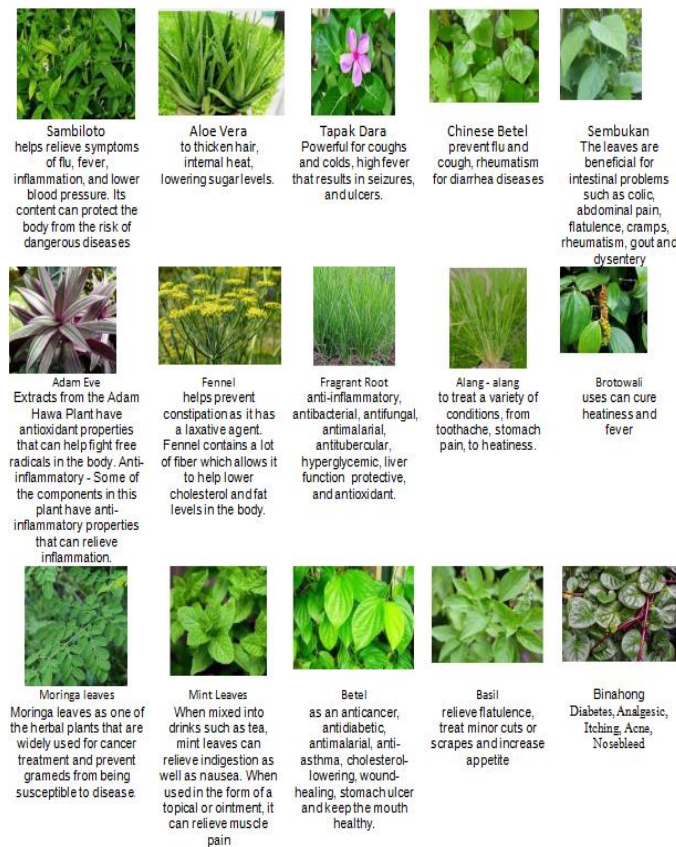


Fig. 3. is a sample of training data used for image retrieval testing on medicinal plant information search.

3.1 Identification Results

Image retrieval results are carried out with variations in the number of clusters and variations in the percentage of feature weights. The test results with 40% weighting variation for shape features and 40% weighting of color features and 20% texture features with a variation of 10 clusters are shown in Figure 4.

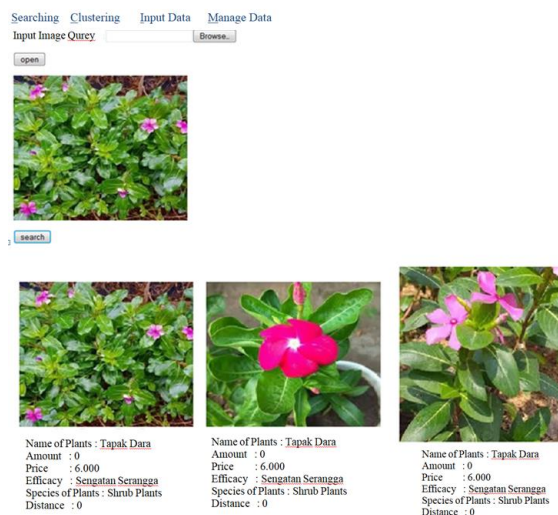


Fig. 4. Medicinal Plant Retrieval Results using Schema W4

Figure 4 display the results of medicinal plant identification using the similarity value with the Euclidean distance method. Retrieval has a similarity ranking of 1 to 10 from the database of medicinal plant leaf images. Testing with variations in weighting and variations in the number of clusters produces different retrieval outputs on the same data.

4.2 Discussion

At this stage, the retrieval accuracy is calculated on the database of medicinal plant leaf images. Furthermore, the retrieval accuracy of the medicinal plant image database is analyzed before and after clustering with a varying number of clusters ranging from 3 to 12 clusters.

3.2.1 Before Clustering

Retrieval accuracy analysis at this stage uses a database of medicinal plant leaf images before clustering with variations in weighting on shape and color features. The accuracy value of retrieval results with the same feature weight and varying feature weights can be seen in Table 1.

Table 1. Percentage of Retrieval Accuracy With Feature Combination On Image Database Before Clustering

Image of Name	Shape Features	Shape and color features	Color and texture features	Shape and texture features	Shape, color and texture features
O1	60,05%	80,20%	70,30%	60,56%	80,50%
O2	60,26%	80,55%	75,10%	70,85%	80,15%
O3	60,02%	80,50%	75,70%	72,33%	81,25%
O4	56,50%	83,30%	76,15%	71,75%	82,65%
O5	59,55%	80,40%	73,47%	72,85%	83,55%
O6	60,35%	80,20%	70,66%	72,32%	83,40%
O7	56,23%	80,60%	75,37%	72,85%	80,65%
O8	60,10%	80,50%	75,80%	71,93%	82,55%
O9	60,25%	80%	75,28%	71,95%	80,03%
O10	60,27%	81,24%	75,10%	71,65%	80,55%

Table 1 shows that the accuracy of image identification with a combination of three features, namely shape, color and texture features, produces the highest average accuracy and the use of one feature produces the lowest image identification accuracy.

3.2.2 After Clustering

The grouping performed on the image is based on similarity, which will narrow down the search space for image data information. The process is very beneficial because it will shorten the identification time and increase the identification accuracy. The level of identification accuracy with the combination of features and clustering can be seen in Table 2.

Table 2. Percentage of Identification Accuracy With Various Combinations Of Image Features With Clustering

Number of clusters	Shape Features	Shape and color features	Color and texture features	Shape and texture features	Shape, color and texture features
3	60,40%	80,10%	72,45%	70,50%	85,60%
4	61,80%	80,25%	72,07%	70,50%	85,40%
5	62,20%	81,35%	73,45%	71 %	85,50%
6	62,50%	81,55%	73,15%	71,30%	85,70%
7	63,80%	82,25%	74,65%	72,40%	85,40%
8	63,30%	82,15%	74,45%	72%	85,80%
9	62,50%	81,46%	74,35%	72,20%	86,80%
10	65,82%	85,55%	77,56%	74,60%	88,60%
11	64,40%	82,15%	76,05%	72,50%	86,50%
12	63,80%	83,46%	74,15%	72%	85,50%

Table 2 shows that the combination of the three features, namely shape, color and texture features as well as variations in the number of clusters of 10 resulted in the highest identification.

accuracy of more than 85%. While the use of only 1 feature at the number of clusters of 3 resulted in the lowest accuracy of less than 62%.

1) The identification accuracy graph on the medicinal plant image database using a combination of shape, color and texture features with a variation in the number of clusters ranging from 3 to 12 is shown in Figure 5.

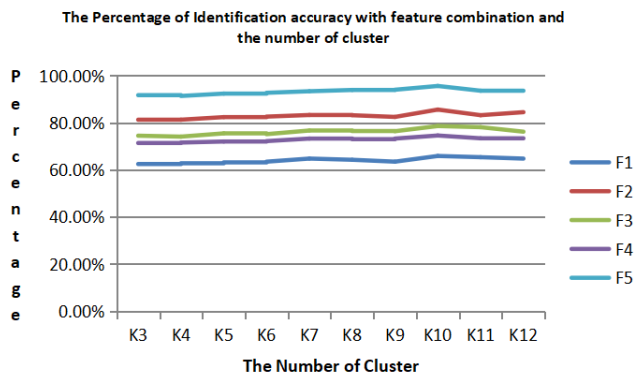


Fig. 5. Graph of Identification Accuracy with the combination and Variation of Number of Clusters

In Figure 5, the identification accuracy tends to increase until it reaches a variation of the number of clusters of 10 and tends to decrease after the number of clusters is greater than 10.

### 3.2.3 Computation Time

The computation time taken to identify an image in the medicinal plant image database is affected by the size of the image database and the database management technique. The average computation time for searching the image database before clustering is with the varying number of clusters shown in Figure 6.

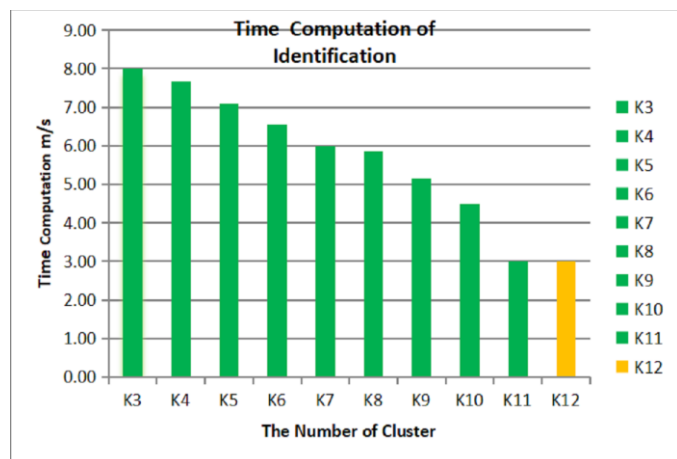


Fig. 6. The Time Computation of Identification

In Figure 6, explains that the more the number of clusters, the faster the computation time. The test results show that identification using three features, namely shape, color and texture features, has a higher accuracy rate than when using

only one or two features. Shape features have the most significant influence on the accuracy of identification results compared to color and texture features. The use of KNN image classification has a significant effect on increasing the speed of the identification process and identification accuracy. Characteristic extraction methods can be developed so as to produce accuracy, precision and recall values with higher percentages.

## 5. DISCUSSION

The accuracy of medicinal plant leaf image retrieval is influenced by several variables including image acquisition, preprocessing, feature extraction methods, feature combinations, feature weighting schemes and clustering techniques in the database. The test results in this study prove that, the use of a combination of features, different feature weighting schemes on each normalized feature value and variations in the number of clusters, will affect the increase in accuracy and speed of the medicinal plant identification process. The results of the tests that have been carried out that the highest level of retrieval accuracy is in the number of clusters 10 with W4 weighting variations, namely 40% shape features and 40% color features and 20% texture features.

This study found that the shape feature in the image of the leaves of medicinal plants has a more dominant factor to determine the level of similarity in the identification process of medicinal plants. As for the color feature as a complementary feature. The average computation time required for retrieval is 5 milli-seconds.

## REFERENCES

- [1] Singh, Krishna, Indra Gupta, and Sangeeta Gupta, (2010). "SVM-BDT PNN and Fourier moment technique for classification of leaf
- [2] Kebapci, Hanife, et all, April (2010). "Plant Image Retrieval Using Color, Shape and Texture Features". The Computer Journal Advance Access published April 9, 2010.
- [3] Gonzales, R. C.; and Woods, R. E. (2008). Digital Image Processing ThirdEdition. Pearson PrenticeHall, New Jersey
- [4] Susilo, A., Web Image Retrieval for Flower Identification with Color Content Grouping, Institut Teknologi Sepuluh November, Surabaya, 2007.
- [5] Vadivel, A.; Majumdar, A.K.; and Shamik, S. (2004). Characteristics Of Weighted Feature Vector In Content-Based Image Retrieval Applications. IEEE.
- [6] Wu, Qingfeng, Changle Zhou, and Chaonan Wang, (2006). "Feature extraction and automatic recognition of plant leaf using artificial neural network". Advances in Artificial Intelligence 3.
- [7] Castleman, K. R. (1996). Digital Image Processing. Prentice Hall Inc. New Jersey.
- [8] Jumi; and Harjoko, A. (2012). Image Similarity Analysis Based on Shape, Color and Texture Feature of Asset Image. International Conference on Computer Science Electronics and Instrumentation, Yogyakarta, Indonesia.
- [9] Xing-yi Huang, Jian Li, and Song Jiang. "Study on identification of rice varieties using computer vision [J]". In: Journal of Jiangsu University (National Science Edition) 2.003 (2004).
- [10] Harpreet Kaur and Baljit Singh. "Classification and grading of rice using multi-class SVM". In: International Journal of Scientific and Research Publications 3.4 (2013), pp. 1-5.
- [11] JD Guzman and EK Peralta. "Classification of Philippine rice grains using machine vision and artificial neural networks." In: World conference on agricultural information and IT, IAALD AFITA WCCA 2008, Tokyo University of Agriculture, Tokyo, Japan, August 24-27, 2008. Tokyo University of Agriculture. 2008, pp. 41-48.

- [12] Zhao-yan Liu, Fang Cheng, Yi-bin Ying, and Xiu-qin Rao. "Identification of rice seed varieties using neural network". In: Journal of Zhejiang University. Science. B 6.11 (2005), p. 1095.
- [13] Chaturika Sewwandi Silva and Upul Sonnadara. "Classification of Rice Grains Using Neural Networks". In: Proceedings of Technical Sessions. Vol. 29. 2013, pp. 9-14.
- [14] Singh, S. K., Bejagam, K. K., An, Y., & Deshmukh, S. A. (2019). Machine-learning based stacked ensemble model for accurate analysis of molecular dynamics simulations. The Journal of Physical Chemistry A. doi:10.1021/acs.jpca.9b03420
- [15] Acharya T, Ray, A. K, 2005, Image Processing Principles and Applications, John Willey & Sons, USA.
- [16] S. G. Wu, F. S. Bao, E. Y. Xu, Y. X. Wang, Y. F. Chang, and Q. L. Xiang, "A Leaf Recognition Algorithm for Plant Classification Using Probabilistic Neural Network," in 2007 IEEE International Symposium on Signal Processing and Information Technology, 2007, pp. 11-16.