

The Grouping Of Facial Images Using Agglomerative Hierarchical Clustering To Improve The CBIR Based Face Recognition System

Muhammad Fachrurrozi

Informatics Engineering Department
Faculty of Computer Science, Universitas Sriwijaya
Palembang, Indonesia
mfachrz@unsri.ac.id

Clara Fin Badillah

Informatics Engineering Department
Faculty of Computer Science, Universitas Sriwijaya
Palembang, Indonesia
09021181320052@students.ilkom.unsri.ac.id

Saparudin

Informatics Engineering Department
Faculty of Computer Science, Universitas Sriwijaya
Palembang, Indonesia
saparudin@unsri.ac.id

Junia Erlina

Informatics Engineering Department
Faculty of Computer Science, Universitas Sriwijaya
Palembang, Indonesia
juniaerlina@gmail.com

Erwin

Computer Engineering Department
Faculty of Computer Science, Universitas Sriwijaya
Palembang, Indonesia
erwin@unsri.ac.id

Mardiana

Law Department
Faculty of Law, Universitas Sriwijaya
Palembang, Indonesia
mardiana_rachman@yahoo.com

Auzan Lazuardi

Informatics Engineering Department
Faculty of Computer Science, Universitas Sriwijaya
Palembang, Indonesia
auzanlazuardi@gmail.com

Abstract— The grouping of face images can be done automatically using the Agglomerative Hierarchical Clustering (AHC) algorithm. The pre-processing performed is feature extraction in getting the face image vector feature. The AHC algorithm performs grouping using linkage average, single, and complete method. Grouping face images can help improve the search speed of the CBIR based face recognition system. The cluster validation test uses the value of Cophenetic Correlation Coefficient (CCC). From the test results, it is known that the complete method has a higher CCC value than other methods, that is equal to 0.904938 with the difference value of 0.127558 on single method and the difference of 0.02291 on the average method. The face recognition system using pre-processing clustering can perform faster face recognition better than without pre-processing clustering.

Keywords— *Clustering; AHC; Single Linkage; Complete Linkage; Average Linkage;*

I. INTRODUCTION

Content Based Image Retrieval (CBIR) is the image process of a database or digital image library in accordance with the visual content of the image [1]. CBIR only focuses on image search using queries on large image databases based on texture, color, shape, and region features. The grouping of facial image can be used to speed up the image search process on facial image recognition system using Image Processing Science. Face recognition is faster and more accurate on CBIR[2] [3][4].

Grouping is divided into two types, namely hierarchy and non hierarchy[5]. Hierarchical Clustering is a grouping algorithm by forming hierarchies of similar data into a tree or dendogram. In Hierarchical Clustering there are two ways of grouping, namely agglomerative and divisive. Agglomerative

clustering process based on the amount of data grouped into hierarchy - hierarchy, then hierarchy - the hierarchy becomes a hierarchical unity.

Extraction facial image feature to get vector feature using Local Binary Pattern (LBP). The distance between the vector features is calculated using the manhattan distance which is subsequently grouped using the Agglomerative Hierarchical Clustering (AHC) algorithm. Manhattan Distance provides relatively higher results than Euclidean Distance with high probability [6].

This study focused on comparing the three methods of the AHC algorithm, namely Single Linkage, Complete Linkage, and Average Linkage in categorizing facial images and helping to improve the speed of face recognition system.

II. RELATED WORK

Grouping using the Hierarchical Clustering algorithm can improve the speed and accuracy of image matching in CBIR [4] that grouped on the semantic image obtained CBIR results increased from 33% to 57% in Fabric image and from 31% to 60% on Sports and Athletic imagery.

Other studies discuss image groupings [7] using K-Means Clustering in 2600 fruit images capable of producing image groupings by combining pixels of the same color for image segmentation to objects effectively and efficiently.

Other studies discussing the multi-object image grouping [8] using the Hierarchical Temporal Memory Network algorithm on 100,000 images with Gaussian noise can increase the speed of 91.4% in identifying the image patterns.

The Content Based Image Retrieval System [3] research uses a combination of K-Means Clustering and Hierarchical Clustering algorithms on large image databases based on

colors, shapes, textures, and patterns. The results show that the combination of K-Means Clustering and Hierarchical Clustering algorithms is able to group images and generate faster searches on large image databases.

III. CLUSTERING

A. Local Binary Pattern (LBP)

LBP represents a pixel which formed by a 3x3 matrix as a comparison between the center pixel and its surrounding pixel which then converted into binary numbers. The comparison assumes that if the surrounded pixel value is greater than the central pixel value than it will be 1 otherwise 0. After we get 8 binary numbers in each pixel then it will be replaced with the decimal form to get the result.

The LBP algorithm formula can be expressed as the following formula:

$$LBP(x_c, y_c) = \sum_{p=0}^7 f(g_p - g_c) 2^p \quad (1)$$

Information:

g_p : central pixel value

g_c : the pixel value around the center

p : number of pixels around the center

And the function $f(x)$ is defined as follows:

$$f(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases}$$

B. Agglomerative Hierarchical Clustering

Agglomerative Hierarchical clustering is a clustering algorithm based on the proximity distance between two images into a hierarchy. This process repeats itself until it gets some hierarchy. The hierarchy with the closest distance is combined into one hierarchy. The proximity to the new hierarchy then recalculated and the closest hierarchy is merged again. The process is repeated until all the data (object) clustered into one hierarchy.

Calculating the spacing between two images using the *Manhattan Distance* formulated in formula (2):

$$d = \sum_{i=1}^n |u_i - v_i| \quad (2)$$

where:

d : the distance between the image of u and v .

n : number of variables.

u_i : the value of u on i variable.

v_i : the value of v on i variable

The distance between images is written into a matrix called distance matrix. In order to determine the distance between the two clusters, Agglomerative Hierarchical clustering has 3 methods of grouping data, namely:

1. Single Linkage

Single Linkage classifies data based on the closest distance (Min) between the hierarchy. Single Linkage can be formulated in formula (3):

$$d(uv)w = \text{Min} [d(uw), d(vw)] \quad (3)$$

where :

u : the image of u

v : the image of v

w : the image of w

$d(uv)w$: the distance between the hierarchy uv and w .

$d(uw)$: the distance between the hierarchy u and w .

$d(vw)$: the distance between the hierarchy v and w .

2. Complete Linkage

Complete Linkage categorizes data by the furthest distance (Max) or the maximum distance between hierarchies. Complete Linkage can be formulated in formula (4):

$$d(uv)w = \text{Max} [d(uw), d(vw)] \quad (4)$$

where :

u : the image of u

v : the image of v

w : the image of w

$d(uv)w$: the distance between the hierarchy uv and w .

$d(uw)$: the distance between the hierarchy u and w .

$d(vw)$: the distance between the hierarchy v and w .

3. Average Linkage

Average Linkage classifies data based on the average distance between the hierarchy. Average Linkage can be formulated in formula (5):

$$d(uv)w = \frac{d(uw) + d(vw)}{2} \quad (5)$$

where :

u : the image of u

v : the image of v

w : the image of w

$d(uv)w$: the distance between the hierarchy uv and w .

$d(uw)$: the distance between the hierarchy u and w .

$d(vw)$: the distance between the hierarchy v and w .

IV. METHODOLOGY

A. Data

The data used as many as 200 images from 20 people, each person taken 10 images with different sides and the same background. With the rules on the image of the face still looks both eyes, nose and mouth. The image is used as training data with dimensions of 150x150 pixels. Examples of face image data can be seen in figure 1.



Fig. 1. Example of image data

B. General System

Figure 2 is a block diagram of a clustering system, where there are 5 stages in the grouping of face images.

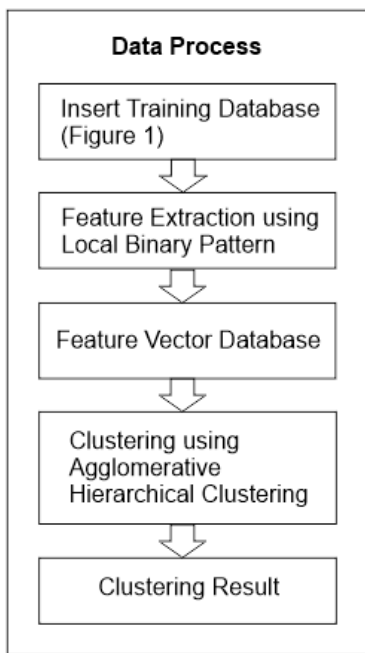


Fig. 2. General System Diagram

In general, the steps of clustering in this research is as follows:

1. Collect face image.
2. The feature extraction process using LBP method to get the characteristic of face image then transformed into vector feature form
3. Vector feature will be stored in the database.
4. Then do the clustering process using the AHC method on the vektor of the face image in the database.
5. Save clustering result.

V. IMPLEMENTATION AND RESULT

A. Clustering

The grouping process starts after getting the value of the vector feature. Grouping begins by calculating the distance between objects using formula (2), to get all the distance between objects can be calculated and then written into a matrix called distance matrix.

Distance matrix is grouped into several groups according to the three methods. Similarities, the closest distance to the formula (3), the farthest distance by the formula (4), or by the mean distance between the image and the other image by the formula (5) The result of the grouping is the grouped vector feature database.

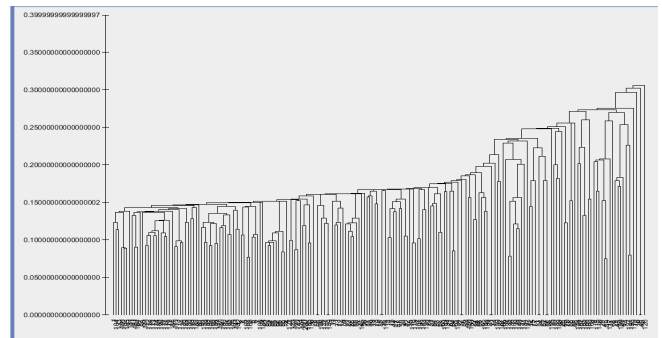


Fig. 3. Result of single linkage

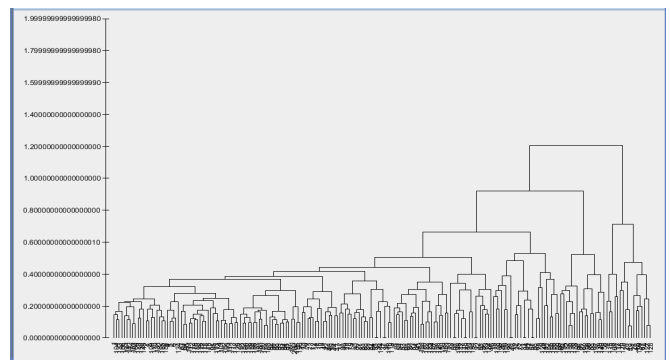


Fig. 4. Result of complete linkage

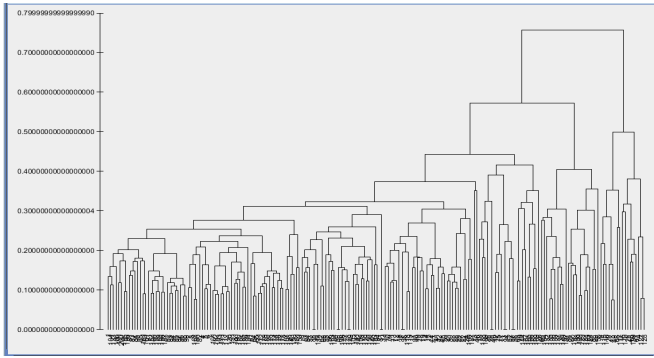


Fig. 5. Result of average linkage

Figure 3, 4, and 5 are the result dendrograms of the face image groupings of the three methods of the AHC algorithm. All three methods produce 199 clusters. The validation of the cluster results of each method using Cophenetic Correlation Coefficient (CCC).

CCC is the correlation coefficient between the original matrix elements of distance matrix and the resulting elements of the dendrogram (Cophenetic matrix)[9]. CCC can be formulated in formula (6):

$$coph = \frac{\sum_{i < k} (d_{ik} - \bar{d})(d_{cik} - \bar{d}_c)}{\sqrt{\left[\sum_{i < k} (d_{ik} - \bar{d})^2 \right] \left[\sum_{i < k} (d_{cik} - \bar{d}_c)^2 \right]}} \quad (6)$$

Information :

$coph$: cophenetic correlation coefficient

d_{ik} : the distance matrix between object i and k

\bar{d} : average of d_{ik}

d_{cik} : distance of cophenetic object i and k

\bar{d}_c : average of d_{cik}

The resulting value of cophenetic correlation coefficient ranges between -1 and 1. The closer to the value of 1 means the quality resulting from the clustering process is said to be good, whereas if the value of CCC approaching value -1 means that the resulting quality of the clustering process is not good.

TABLE I. COPHENETIC CORRELATION COEFFICIENT VALUES

Single Linkage	Complete Linkage	Average Linkage
0.777381	0.904938	0.882028

In Table I the value of *Cophenetic Correlation Coefficient* obtained is single linkage is 0.777381, *complete circle* is 0.904938, and the *average linkage* is 0.882028. All three methods have a CCC value close to 1, which means the three methods are able to group well on the research data (200 face images). However, from the three methods it is known that the *Complete Linkage* method has a higher CCC value

than the other two methods with a difference of value of 0.127558 on *single linkage* method and the difference of 0.02291 on the *average linkage* method. This shows the complete linkage method is better in grouping research data of 12.76% of the *single linkage* method and by 2.29% of the *Average Linkage* method.

B. Implementasi Pra-Pengolahan Clustering pada Sistem Pengenalan Citra Wajah

General system of face recognition can be seen in figure 6.

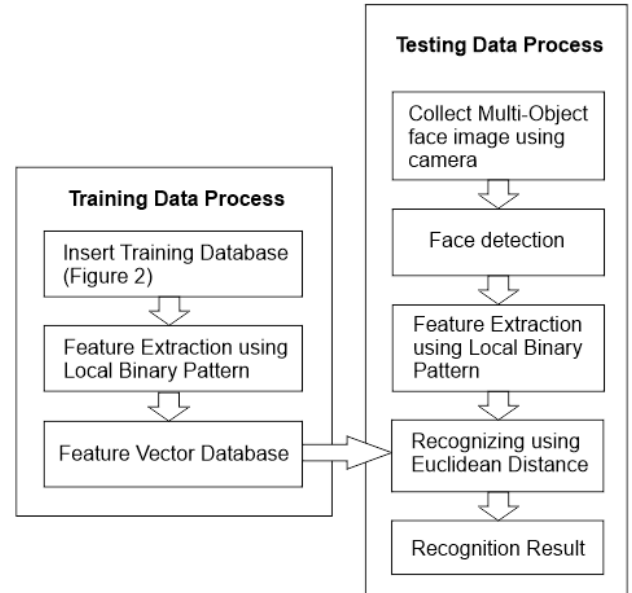


Fig. 6. General System of Face Recognition

In general, the steps of data retrieval in this research is as follows:

1. Collect face image.
2. The feature extraction process using LBP method to get the characteristic of face image then transformed into vector feature form which will be stored in the database.
3. Then do the clustering process using the AHC method on the vector of the face image in the database. Clustering results that have been obtained later used as a comparator value of calculating the distance for face recognition.

In the face recognition phase the process of the system is as follows:

1. Open webcam to detect faces. The process of face detection and recognition is done in real-time.
2. The detected face then captured.
3. The feature extraction process is done using Local Binary Pattern method to get the feature of the face image then transformed into the vector feature form.

- Then, do the process of recognition by calculating the distance between the new face image features and features of the existing on the database by using Euclidian distance which then matched with the clustering results.

Face recognition typically performs face recognition by matching the vector feature of test data with existing vector features from training data. The search process takes quite a long time, because the search is done on all vector features of the training image. Meanwhile, facial image recognition using pre-processing process of AHC clustering performs the search process of vector feature by searching for cluster approaching vector feature of training data, so the search is not done entirely on train data. This further saves the computing time of the system doing face recognition.

The process of recognizing facial image using pre-processing clustering can be seen in figure 7.

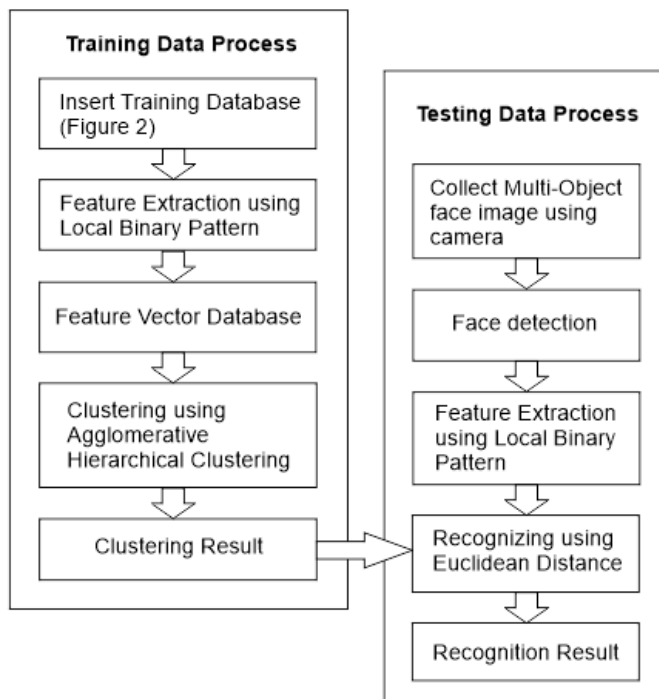


Fig. 7. Clustering in Face Recognition System

The CBIR speed time test using AHC is calculated using the start time and stop time of program execution. Computation time is done to get the computation time of each method process in assisting face recognition system in recognizing face image. CBIR speed computation time test without clustering and CBIR with clustering using AHC is done with the formula in equation (7), that is:

$$R_{time} = \frac{\sum W}{n} \quad (7)$$

Keterangan :

R_{time} : the average of time

n : number of variables

W : computing's time

Tables II and III are the results of calculating the computational time of the face recognition system using clustering and without clustering. The average value of time calculation result is calculated using formula 7. In Table II facial recognition experiments on 1 object (1 face in 1 image) the same without using clustering process as much as 10 times of facial recognition experiments and using clustering process 10 times of facial recognition experiments on each method. While in table III facial recognition experiments on 6 objects (6 face in 1 images) the same without using the clustering process as much as 10 times of facial recognition experiments and using clustering process 10 times of facial recognition experiments on each method. Testing the computing time of this facial recognition system proves whether the three proposed methods can help reduce the time of face recognition system in recognizing face image.

TABLE II. TABLE II. FACE RECOGNITION IMAGE USING 1 OBJECT

Recognition n to -	Computation time without clustering process (s)	Computation time uses clustering process (s)		
		Single Linkage	Complete Linkage	Average Linkage
1	0.781	0.759	0.607	0.779
2	0.781	0.764	0.754	0.765
3	0.797	0.754	0.075	0.766
4	0.829	0.755	0.754	0.766
5	0.797	0.749	0.074	0.765
6	0.766	0.764	0.772	0.750
7	0.766	0.771	0.751	0.765
8	0.766	0.779	0.754	0.765
9	0.781	0.754	0.750	0.765
10	0.813	0.753	0.764	0.766
Average	0.7877	0.7602	0.6055	0.7652

TABLE III. TABLE II. FACE RECOGNITION IMAGE USING 6 OBJECT

Recognition n to -	Computation time without clustering (s)	Computational time using clustering process (s)		
		Single Linkage	Complete Linkage	Average Linkage
1	0.779	0.722	0.766	0.757
2	0.897	0.751	0.750	0.758
3	0.766	0.765	0.750	0.772
4	0.875	0.754	0.766	0.755
5	0.781	0.757	0.765	0.757
6	0.766	0.755	0.751	0.758

7	0.780	0.753	0.797	0.752
8	0.766	0.776	0.750	0.756
9	0.782	0.751	0.750	0.776
10	0.766	0.770	0.750	0.753
Average	0.7958	0.7554	0.7595	0.7594

In table II test for face recognition not real-time on 1 object, it is known that the three methods are able to increase the speed of face recognition system with single linkage method by 3.49%, complete linkage method is 23.13%, and the average circumference method is 2.25%.

In testing table III for face recognition not real-time on 6 objects, it is known that the three methods are able to increase the speed of face recognition system with single linkage method 5.08%, complete linkage method is 4.56%, and the average linkage method is 4.57%.

VI. DISCUSSION

In this study it is expected that pre-processing of clustering using AHC on facial recognition system can increase the computation time speed of facial recognition system. However, in the results of this study found that the process of pre-processing clustering using AHC can only slightly increase the speed of computing time recognition system. Results from the cluster validity test, the three methods on the AHC algorithm are able to group the data well. The computing speed is slightly influenced by the methods used in similar image searches on facial recognition systems.

VII. CONCLUSION

From the results of testing new AHC method with single method, complete, and average of good association of facial image grouping. Among the three methods, it is known that the Complete Linkage method is better than the other two methods (single and average linkage).

The face recognition system using pre-processing clustering can perform faster face recognition better than without pre-processing clustering.

REFERENCE

- [1] R. Chaudhari and A. M. Patil, "Content Based Image Retrieval Using Color and Shape Features," *Int. J. Adv. Res. Electr. Electron. Instrum. Eng.*, vol. 1, no. 5, pp. 386–392, 2012.
- [2] A. K. V. M. N. Anbazhagan, "Image Clustering and Retrieval using Image Mining Techniques," 2010.
- [3] V. S. V. S. Murthy, E. Vamsidhar, J. N. V. R. S. Kumar, and P. Sankara Rao, "Content Based Image Retrieval using Hierarchical and K-Means Clustering Techniques," *Int. J. Eng. Sci. Technol.*, vol. 2, no. 3, pp. 209–212, 2010.
- [4] S. Pandey, P. Khanna, and H. Yokota, "Clustering of hierarchical image database to reduce inter-and intra-semantic gaps in visual space for finding specific image semantics," *J. Vis. Commun. Image Represent.*, vol. 38, pp. 704–720, 2016.

- [5] R. Nainggolan, "Algoritma Modified K-Means Clustering pada penentuan Cluster Centre Berbasis Sum of Squared Error (SSE)," 2014.
- [6] C. C. Aggarwal, A. Hinneburg, and D. A. Keim, "On the surprising behavior of distance metrics in high dimensional space," *Database Theory – ICDT 2001*, pp. 420–434, 2001.
- [7] R. Mente, B. V. Dhandra, and G. Mukarambi, "Color Image Segmentation and Recognition based on Shape and Color Features," *Int. J. Comput. Sci. Eng.*, vol. 3, no. 1, pp. 51–56, 2014.
- [8] L. Xie, K. Yang, and X. Gao, "Multi-object recognition by optimized hierarchical temporal memory network," *Opt. - Int. J. Light Electron Opt.*, vol. 127, no. 19, pp. 7594–7601, 2016.
- [9] A. Rodrigo and C. Tadeu, "A cophenetic correlation coefficient for Tocher's method," no. 1, pp. 589–596, 2013.