



Building Synonym Set for Indonesian WordNet using Commutative Method and Hierarchical Clustering

Valentino Rossi Fierdaus¹, Moch. Arif Bijaksana², Widi Astuti³

Faculty of Informatics, Bachelor of Informatics Engineering, Telkom University, Bandung, Indonesia

Email: ¹ valentinorfs@student.telkomuniversity.ac.id, ² arifbijaksana@telkomuniversity.ac.id,

³ astutiwidi@telkomuniversity.ac.id

Email Coresspondence: valentinorfs@student.telkomuniversity.ac.id

Abstract—WordNet is a compilation of Synonyms Set (synset), which consists of the words that have the same synonymous. The development of Indonesian WordNet has a goal to build an application that can accommodate and exhibit the relation of words. Synonym Set is a set composed of one or more words that have a similar meaning or synonym relation originated from the Indonesian Thesaurus. In previous studies, the establishment of synsets were transmitted with several approaches, one of which was the cluster ring to produce synsets and WSD (Word Sense Disambiguation). In this research, research is held up to discover the semantic similarities between words in the Indonesian Thesaurus automatically, and also to know the performance of the Agglomerative Hierarchical Clustering method for the development of Indonesian synsets. To calculate performance and evaluation, this research is using the F-measure method involving the gold standard.

Keywords: WordNet, Synset, Indonesian Thesaurus, Agglomerative Hierarchical Clustering, F-Measure.

1. INTRODUCTION

WordNet is a set of several synonyms called Synonym Set (Synset) consisting of words that have equivalent meanings or sense which are interrelated [1]. At first, WordNet is a semantic dictionary that made in English version which was first built by Princeton University, then along with development of technology, WordNet at present is one of the most widely used sources of referral information. Language dictionaries around the world in general is a dictionary that has a focus words while WordNet focuses on the meaning of words or synonyms. In WordNet, several classes of words such as nouns, verbs, adjectives, and adverbs is grouped into a synsets. Lines of words in WordNet can symbolize a meaning which is called synset.

In the process of building WordNet, the first thing to do is produce a synset or collection of synonym that have same meanings [1], that means the words are grouped into a synset according to their meaning. That is because synset is a basic concept that supports the formation of semantic relations in the lexical database [2]. Monolingual resource that used as a lexical resource is Thesaurus, because Thesaurus contains words that have an interrelated synonym relation [1]. Thesaurus that has been through the extraction process, will produce one or more synset. To combine the synset that produced from the previous process, one way to produce the best synset is using the clustering techniques. Therefore, need some further research to find out the performance of clustering techniques in the development of synset for Indonesian WordNet.

Previously, there was a development of Indonesian WordNet using Hierarchical Clustering. In that study, the data that used as input is a synset that was generated from the commutative process, then that data (synset) will be grouped and combined in the Clustering process. However, the data used as input are data generated from the results of manual commutative process. This research will focus on two main stages, the first is the stage to doing synset extraction, and the second stages is the process of combining synsets using clustering technique if in the first stages there is a word produces more than one synset. In the synset extraction process, to produce a valid synset value will use Commutative method using available monolingual resources that is Thesaurus Bahasa Indonesia, this means that if a word k1 has a synonym k2, then k2 must also be a synonym of k1. In fact, commutative relations like this do not always occur in Thesaurus Bahasa Indonesia [1]. And for the second stages, clustering technique that used in this research is Agglomerative Hierarchical Clustering.

The purpose of this research is to find out the semantic similarities between words in Thesaurus Bahasa Indonesia automatically, and also implement the Agglomerative Hierarchical Clustering method on the system to be built to determine the performance of that clustering techniques in the development of synset for Indonesian WordNet.

2. RESEARCH METHODOLOGY

2.1 Research Stages

Flowchart in the process of building synonym set for Indonesian WordNet can be seen at Figure 1.

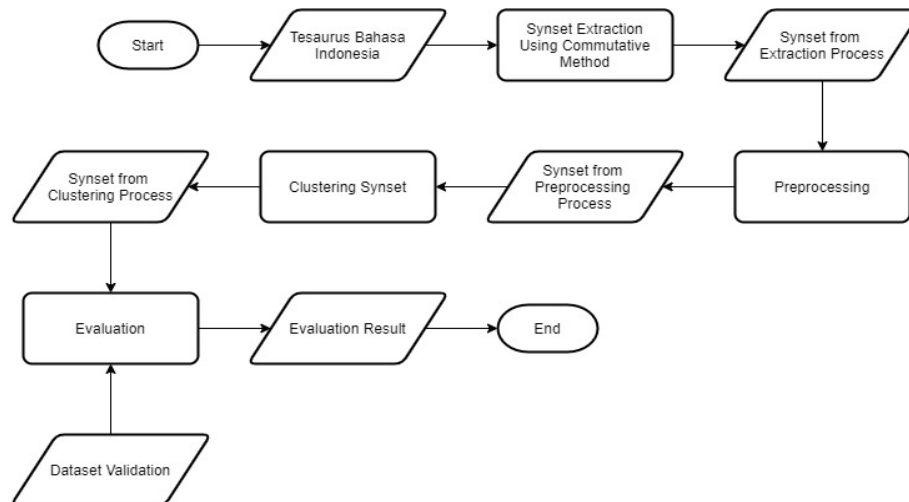


Figure 1. Flowchart System

The flow of the system to be built is depicted at the Figure 1, explained that the initial input is a word from Thesaurus Bahasa Indonesia that would be processed to find out the equivalent meaning with the other words. After that is doing the process of synset extraction that is identification process of the data test using commutative method to produce a synsets. Then, synset that has been through the extraction process will be processed in preprocessing, this process aims to remove excessive characters and spaces in the synset that produced from previous process. Other than that, preprocessing can also make the system read the data properly. After that, the synset that produced from preprocessing will be grouped and combined using Agglomerative Hierarchical Clustering.

The system to be built are expected can be produce the synset from the words that entered as an input, also to find out the semantic similarities between words in Thesaurus Bahasa Indonesia, to find out the performance of clustering techniques in the development of synset for Indonesian WordNet.

2.2 WordNet

WordNet is an online lexical database. This development is based on the theory psycholinguistic human lexical memory. In WordNet, verbs, nouns, adjectives, and adverbs are grouped into a collection of cognitive synonyms (synset), in order to represent different concepts [3].

In research [4], WordNet contains information about 155,000 from various classes of words such as nouns, verbs, adjectives, and adverbs, these words are grouped into a synset according to their meaning. At present, Indonesian WordNet has 1203 (synset) and 1659 unique words in it. The number of semantic relations that can be made from existing synset reaches 2261 relations [5].

2.3 Automatic WordNet Development

In the development of Indonesian WordNet, previously there was a development of Indonesian WordNet based on Linked Data, while the stages of development included identification of data sources, data extraction, data transformation, loading data into relational databases, and mapping relational databases to RDF models [6].

Other than that, previously there was a development of Indonesian WordNet using Hierarchical Clustering. In that study, the data that used as input is a synset that was generated from the manual commutative process, then that data (synset) will be grouped and combined in the Clustering process [7]. While in this study, the data used as input is a word derived from the Thesaurus Bahasa Indonesia which will then be processed using the automatic commutative process. Then, the data (synset) that generated from commutative process will be used to be grouped and combined in the Clustering process.

2.4 Synonym Set (Synset)

Synonym Set (Synset) is a collection that composed of one or more words that have a relation of synonym. Each member of this set can replace each other without changing the sense or meaning of the sentence that contains it [1]. The role of synset becomes very important in WordNet of any language, because all semantic relations connecting a synset, not the word.

[akomodasi, fasilitas]

Figure 2. Example of synonym sets (synset)

Figure 2 is an example of a word that has a synonym relation that is the word “akomodasi” and “fasilitas”. These two words have the same meaning which is something that serves to meet the needs for launch the process of a



particular business. Synonym sets not always have more than one member, synonym sets also have a single member, this can happen because the entry word has no synonym.

2.5 Commutative Method

Synset is valid if there has commutative concept. In Princeton WordNet, synonym relations should be commutative which means that if a word *k1* has a synonym *k2*, then *k2* also must be a synonym of *k1* [1]. Finding a synset that has a valid value is done using a matrix table. Table 1 is shows the matrix table of the word “Fana”. Sense of “fana” is “sementara”, and “kontemporer”, which that words will be used as rows and columns [8].

Table 1. Matrix table “fana”

	fana	sementara	kontemporer
fana	T	T	T
sementara	T	T	T
kontemporer	F	T	T

Explanation of Table 1, the cells that has a true value (T) indicating that the relation of the words contained in the row and column is commutative, and for the cells that has a false value (F) they have the opposite meaning. For example, the row of “temporer” and the column of “fana” has the cell that valued false (F), this means “fana” does not the sense of “temporer”. From that matrix table produced a synset that is [fana, sementara].

2.4 Synset Extraction

Before the preprocessing stage is carried out, there is an extraction process using the commutative concept. The extraction process is carried out in several stages of the algorithm as follows [9].

- a. Searching for a sense of the entry word.
The selected word is then searched for its meaning using a thesaurus. The chosen word for example is “ahad” which has a sense of “minggu”, “esa”, “satu” and “tunggal”.
- b. Searching for synonyms on every sense from point (a).
Every sense in point 1 will be searched for its synonyms.
- c. Searching for “ahad” in the sense that being sought.
The two Senses of the word "ahad" are then searched for a set of synonyms. As seen in Table 2.

Table 2. The synonym of “ahad”

Entry	Synonym
ahad	minggu esa, tunggal, satu
minggu	ahad, pekan
esa	ahad, satu, tunggal
tunggal	ahad, esa, satu
satu	ahad, eka, esa, homo-, iso-, mono-, se-, suatu, tunggal, uni, unik, wahid

The next step is to look for synonyms that have the keywords: “ahad”, i.e. “minggu”, “esa”, “satu”, and “tunggal”, then these words and “ahad” themselves will be made into columns and rows in the matrix.

- d. Identify the prospective synset to be sought.
At this stage the process is done is to look for candidates for the synset that can be generated from each item from the words in the dataset.
 - First sense
 - “ahad”, “minggu”
 - Second sense
 - “ahad”, “esa”
 - “ahad”, “satu”
 - “ahad”, “tunggal”
 - “ahad”, “esa”, “satu”
 - “ahad”, “esa”, “tunggal”
 - “ahad”, “satu”, “tunggal”
 - “ahad”, “esa”, “satu”, “tunggal”
- e. Determine whether every word in the prospective synset has a commutative relationship.
At this stage the process is carried out is to determine the prospective synset that has a commutative relationship. The commutative relationship applies to synset candidates who have more than one set.



- f. Elimination of candidate synset which is a subset of the other synset.
At this stage the process carried out is the elimination of the candidate synset which is a subset of the other synset.
- g. Take the remainder of the elimination synset candidate.
In this process the remaining synset candidates will be taken and used as a synset. For example in "ahad" the resulting final synset is . .
- First sense
 - "ahad", "minggu"
[ahad, minggu]
 - Second sense
 - "ahad", "esa", "satu", "tunggal"
[ahad, esa, satu, tunggal]

2.5 Hierarchical Clustering

Hierarchical clustering is some method of clustering that aims to grouping some data based on the concept of hierarchy. In this method, the two closest groups will be combined in each iteration [10]. In this research, Agglomerative Hierarchical Clustering will be used as Hierarchical Clustering.

2.5 Agglomerative Hierarchical Clustering

Agglomerative Hierarchical Clustering is an bottom-up approach [11], that means where clusters have sub-clusters, which in turn have sub-clusters, etc.. However, a modification needs to be done, because the goal is not to cluster the candidates into a single cluster, in which data will be grouped based on distance value and the clustering process will be stopped after it reached a condition decided by threshold value [2] [12]. Calculation of the distance values can be seen in Equation 1.

$$DistanceValue = \frac{SimilarWords}{UniqueWords} \quad (1)$$

Unique Words are the number of words that exist in the two synonym sets that are compared, and Similar Words are the same number of words in the two synonym sets that are compared. The clustering process will stop if the threshold value is greater or equal to than the maximum distance value. Calculation of the threshold value can be seen in Equation 2.

$$Threshold = coefficient \times firstmaxdistancevalue \quad (2)$$

In Equation 2, threshold value is obtained by multiplying the coefficient value with the maximum distance value that obtained from the first iteration. And each coefficient value can be changed manually in the range of 0.1 to 1.0.

This Clustering process needs to be done to grouping the synset on redundant dataset to produce a better synset. The purpose of this clustering is to group based on the largest similarity value and the largest distance value. Similarity value is obtained from the same number of words between one synset with another synset. Indexes indexes that have the same distance value will be combined with other indexes that have the largest distance value. The clustering process will stop if the threshold value is greater or equal to than the maximum distance value [2]. Distance value and threshold values are generated from Equation 1 and Equation 2. The algorithm used in this study can also be described in the form of a flowchart depicted in the Figure 3 [7].

```

1. Data: Data Set
2. Result: Synset
3. synset 1; synset 2;
4. similarity = 0;
5. unique word = 0;
6. coefficient = [0.1-1.0];
7. while data set do
8.   if synset 1 == synset 2 then
9.     similarity = similarity + 1;
10.  unique word = unique word + 1;
11.  distance value max = similarity / unique word;
12.  threshold = coefficient x distance value max;
13.  while distance value max > threshold do
14.    if distance value max then
15.      synset 1 join synset 2;
16.    Else
17.      distance value max;
18.  Else
19.    Find another synset;

```

Figure 3. Agglomerative Hierarchical Clustering Algorithm



2.5 Tesaurus Bahasa Indonesia

Thesaurus contains a set of words that are related together. Basically, a thesaurus is a means to divert ideas into words, or vice versa. Thesaurus is distinguished from the dictionary. In the dictionaries can be found information about the meaning of words, while in the thesaurus words can be used to express the ideas of the users. Therefore, the thesaurus can help the users to express the ideas according to what is meant [13]. In this research, the thesaurus used is the Tesaurus Bahasa Indonesia in pdf format which was published in 2008.

2.6 Data Test

In the testing phase, 80 words will be used as the data test that taken randomly from the thesaurus, these words will be extracted to produce a valid synset. The selected word is then processed by the system and will produce one or more synset. The list of words used as data test is shown in Table 3.

Table 3. Words for the data test

abah	bahu	cahar	fana	harap	kafilah	nahas	sumpah
aborsi	bajul	ceker	fiber	harfiah	kembang	nasrani	tenda
abrasi	bantahan	dam	fiksi	ialah	lamaran	oral	tending
acuh	bentak	dandang	gagu	idealis	langgar	orasi	tikar
aduk	binatang	darat	gajih	istana	lompat	pacul	tukar
agar	bingkai	edan	galat	jahiliah	madukara	padas	tukar
ahad	bisik	eka	gas	jalin	maharani	pahala	tulis
ampun	buat	eks	gelar	jambe	mambang	pukul	tusuk
arwah	cabuk	eksponen	gosok	kabin	mampir	qadim	umpat
ayam	cadar	faksi	handuk	kabung	nigari	sulam	usap

The data test will be extracted using commutative methods to produce one or more synset which will be processed using Agglomerative Hierarchical Clustering.

2.7 Gold Standard

Gold Standard aims to find out how much the correlation between the score issued by the system and the relevance of the words being tested. The gold standard value is obtained from a collection of human opinions. This value is used as a reference measurement of similarity between words. In this study, the gold standard used is the result of validating synonym sets performed by lexical experts (lexicographers). The validation is done very carefully so that it can be used as a comparison for the results of the system as a measure of accuracy [12].

2.8 F-Measure

F-measure is a popular performance metric, especially for tasks with unbalanced data sets [14]. This F-measure method involves the precision method and the recall method. For the calculation of precision (P) and recall (R) can be seen in equations 3 and 4 [15].

$$P = \frac{\text{jumlahbenardariprediksi}}{\text{jumlahyangdiprediksi}} \quad (3)$$

$$R = \frac{\text{jumlahbenardariprediksi}}{\text{jumlahseharusnya}} \quad (4)$$

The F-Measure method calculates multiple propositions multiplied by the results of the first method (precision) and the second method (recall) divided by the sum of that two. The calculation of F-Measure can be seen in equation 5.

$$F - Measure = 2 \times \left(\frac{P \times R}{P + R} \right) \quad (5)$$

Precision is taken from the calculation of the number of correct words in the synset that has been produced by the system adjusted to the gold standard generated from manual calculations and then divided by the number of words in the resulting synset, while the recall is taken from the calculation of the number of correct words in the synset that has been generated by the system adjusted to the results of the manual calculation then divided by the number of words in the synset that has been calculated manually by humans or the gold standard.

3. RESULT AND DISCUSSION

In this section, will explain the related things about the testing result and the analysis of testing result that have been carried out.



3.1 Testing Result

The method used for the evaluation process is using the F-measure method. In addition, testing is done by changing the coefficient value from the range 0.1 to 1.0. This test is carried out to find out which coefficient values can combine a set of synonyms exactly according to their sense. The test results that have been carried out by changing the coefficient values from the range 0.1 to 1.0 produce data that can be seen in Table 4.

Table 4. Comparison of Coefficient Values

Koefisien	Max Similarity	Max Distance Value	Number of Loop	Number of Synsets
0.1	0	0	34	79
0.2	1.0	0.11	33	80
0.3	1.0	0.17	31	82
0.4	1.0	0.2	28	85
0.5	1.0	0.25	22	91
0.6	1.0	0.33	13	100
0.7	2.0	0.4	11	102
0.8	3.0	0.43	10	103
0.9	3.0	0.5	2	111
1.0	3.0	0.5	2	111

Based on the data in Table 4, it appears that the coefficient values 0.9 and 1.0 have the same number of synset and number of loops, this is because the same maximum distance value. The clustering process continues because the coefficient value is still lower than the maximum distance value generated from each loop. As explained earlier, the clustering process will stop if the threshold value is greater or equal to than the maximum distance value.

Other than that, clustering performance testing process has been carried out on commutative and clustering data set results with a range of coefficient values 0.1 to 1.0. Here is a comparison of performance data generated from the tests that have been carried out.

Table 5. The testing results using F-Measure

Data Test	Precision (%)	Recall (%)	F-Measure (%)
Before clustering	51.3	62.11	56.19
Coefficient 0.1	84.81	72.83	78.36
Coefficient 0.2	83.75	72.83	77.91
Coefficient 0.3	85.37	76.09	80.46
Coefficient 0.4	84.71	78.26	81.36
Coefficient 0.5	81.32	80.43	80.87
Coefficient 0.6	80.0	86.96	83.33
Coefficient 0.7	78.43	86.96	82.47
Coefficient 0.8	78.64	88.04	83.08
Coefficient 0.9	77.48	93.48	84.73
Coefficient 1.0	77.48	93.48	84.73

Table 5 is a comparison of performance data generated from the tests that have been carried out. The clustering performance testing process has been carried out on commutative and clustering data set results with a range of coefficient values 0.1 to 1.0. This is done to find out how the performance of each process by calculating the synset generated by the system.

Table 6. Comparison of synset

Word	Synset Validation	Synset Commutative	Synset Clustering
agar	[agar, biar, supaya]	[agar, biar, supaya] [agar, mudah-mudahan]	[agar, biar, supaya] [agar, mudah-mudahan]
jalin	[jalin, anyam, kepang, rangkai] [jalin, lilit]	[anyam, jalin, rangkai]	[jalin, anyam, rangka]
tiup	[tiup, hembus]	[hembus, tiup] [sembur, tiup]	[tiup, hembus] [tiup, sembur]

Table 6 is a comparison between Synset Validation, Synset Commutative, and Synset Clustering. Synset Validation synset that used as validation, Synset Commutative is a synset that generated from automatic commutative process, and Synset Clustering is a synset that generated from clustering process.

3.2 Analysis of Testing Result



Referring to the data shown in Table 4, the results of testing on each coefficient value produce different values in each aspect, for example the difference in the amount of synset and Distance Value always increases along with the coefficient value which also increases. Inversely, the number of loops in the clustering process always decreases even though the coefficient value continues to increase, this happens because the clustering process will continue as long as the coefficient value is less than the maximum distance value that generated from each loop.

Other than that, the performance in the data set from the smallest coefficient to the largest always increases and reaches the optimum point at Coefficients 0.9 and 1.0, because on these coefficients the value of Max Distance Value, Number of Loops, and Number of Synset has been produced does not change.

Based on the data listed in Table 5, building synonym set using commutative method (before the clustering process) with predetermined data test, obtained a value of Precision at 51.3%, Recall at 63.11%, and F-Measure at 56.19%, where that value has a big difference compared to the performance possessed by the clustering process. For example, the value that generated from the optimum coefficient has a value of Precision at 77.48%, Recall at 93.48%, and F-Measure at 84.73%. This could be indicates that the clustering process is very useful for the process of combining synset to be better.

Then based on Table 6, 3 samples were taken that were used for comparison between validation data, data from commutative process, and data from clustering process. That table shows a different similarity of generated synset. It certainly affects the evaluation value, it means that if the words in validation synset and the synset that produced by the system have significant differences, the resulting evaluation value will be lower. However, if the validation synset and the synset that produced by the system have a similarity, then the evaluation value that will be generated can be even greater.

4. CONCLUSION

Threshold value can be found from the coefficient. Threshold value used to stop looping in the clustering process if the threshold value is greater or equal to than the maximum distance value. Based on data obtained from tests that have been done, the clustering process is very useful for the process of combining several synset to be better synset, as evidenced by the F-Measure value that produced from the clustering process is at 84.21% compared to the F-Measure value produced from the process that was passed without clustering which was only at 56.19%. In addition, the coefficient values in the range 0.1 to 1.0 in the clustering process obtained the coefficient value 0.9 as the most optimal coefficient because of the things that have been explained in the previous chapter. However, the optimum coefficient value can vary depending on the dataset used in the testing.

The suggestion for further research, it is expected to do the addition of a list of words to be used, this has a purpose in order to determine the optimum coefficient value and also to measure the performance of Agglomerative Hierarchical Clustering on a bigger data scale. And then, further research is expected to be used another clustering method in the development of synset for Indonesian WordNet.

REFERENCES

- [1] Gunawan, "Akuisisi Gloss Berbasis Ekstraksi Synonym Set Menggunakan Supervised Learning," Institut Teknologi Sepuluh November, 2016.
- [2] A. Saputra and others, "Building synsets for Indonesian Wordnet with monolingual lexical resources," in *2010 International Conference on Asian Language Processing*, 2010, pp. 297–300.
- [3] G. A. Miller, "WordNet: a lexical database for English," *Commun. ACM*, vol. 38, no. 11, pp. 39–41, 1995.
- [4] C. Fellbaum, "WordNet," in *Theory and applications of ontology: computer applications*, Springer, 2010, pp. 231–243.
- [5] U. Indonesia, "WordNet Bahasa Indonesia," 2008. <http://bahasa.cs.ui.ac.id/> (accessed Jul. 26, 2019).
- [6] H. Hendrik and A. B. Cahyono, "Model WordNet Bahasa Indonesia berbasis Linked Data," *J. Nas. Tek. Elektro dan Teknol. Inf.*, vol. 6, no. 1, pp. 8–14, 2017.
- [7] J. Priyatno, "Clustering Synonym Sets in English WordNet," Universitas Telkom, 2018.
- [8] D. J. Restina, "Pembangunan Synonym Set untuk WordNet Bahasa Indonesia dengan Menggunakan Metode Komutatif," *Indo-JC*, vol. 4, no. 2, 2019.
- [9] I. P. P. Ananda, "Pembangunan Synsets untuk WordNet Bahasa Indonesia dengan Metode Komutatif," Universitas Telkom, 2018.
- [10] K. Sasirekha and P. Baby, "Agglomerative hierarchical clustering algorithm-a," *Int. J. Sci. Res. Publ.*, vol. 83, p. 83, 2013.
- [11] D. Müllner, "Modern hierarchical, agglomerative clustering algorithms," no. 1973, pp. 1–29, 2011.
- [12] L. D. Anggaraini, "Analisis Pembangunan Word Sense pada WordNet Bahasa Indonesia Menggunakan Metode Hierarchical Clustering," Bandung, 2019.
- [13] T. Redaksi, "Tesaurus Bahasa Indonesia Pusat Bahasa," *Pus. Bahasa, Dep. Pendidik. Nas.*, 2008.
- [14] Y. Nan, K. M. Chai, W. S. Lee, and H. L. Chieu, "Optimizing F-measure: A tale of two approaches," *arXiv Prepr. arXiv1206.4625*, 2012.
- [15] D. R. Musicant, V. Kumar, A. Ozgur, and others, "Optimizing F-Measure with Support Vector Machines.," in *FLAIRS conference*, 2003, pp. 356–360.