

# MALAY PART-OF-SPEECH TAGGING: AN ME-BASED APPROACH

Juhaida Abu Bakar<sup>a</sup>, Khairuddin Omar<sup>b</sup>, Mohammad Faizul Nasrudin<sup>b</sup>, Mohd Zamri Murah<sup>b</sup>

<sup>a</sup>School of Computing, College of Arts and Sciences, UUM, 06010, Sintok, Kedah, Malaysia

<sup>b</sup>Center for Artificial Intelligence Technology, Faculty of Information Science and Technology, UKM, 43600, Bangi, Selangor, Malaysia

\*Corresponding author  
juhaida.ab@uum.edu.my

## Abstract

Research on Malay Part-of-Speech (POS) tagging has greatly increased over the past few years. Based on previous literature, POS-tags are known as the first phase in the automated text analysis; and the development of language technologies can barely initiate without this initial phase. Malay language can be written in either the Roman or Jawi scripts. We highlight the existing POS-tags approaches and techniques; and suggest the development of Malay Jawi POS-tags using ME-based approach – using specific contextual information of Malay corpora that has been written in Jawi script. We conduct our test on NUWT Corpus. It has been found out that the ME-based approach reaches an accuracy level of 89.30% in average; and yields the precision and recall rates of 94% for the highest level of accuracy achieved.

Keywords: NLP pipeline task, POS-tags, tagging approach, Malay language, Jawi

## 1.0 INTRODUCTION

Assigning syntactic categories to words is an important pre-processing step for most Natural Language Processing (NLP) applications [1]. Part-of-speech tagging or POS-tags is an important feature in NLP for word-category analysis. Effective analysis of Malay corpora can thus, be maximized through POS-tags; regardless of the writing system – the Roman (Rumi) or the Jawi script. It is generally accepted that the application of POS-tags in NLP applications can greatly improve the quality of NLP tasks. That being said, developing high-quality and fast tagging systems is still deemed to be a problem; despite the applications of several different POS-tags models and methods in various languages.

POS-tagging is the process of contextually assigning syntactic categories (noun, verb, etc.) with the most probable sequence to each word in a sentence. This task is a complex algorithmic process since one particular word might be associated with several possible tags. For example, the Malay word, “*menggembirakan*” (gloss: delightful) can be a verb (as in “*Sara menjalani kehidupan yang menggembirakan di China*”) (gloss: Sara lived a happy life in China) or an adjective (as in “*Kejayaan Lim sungguh menggembirakan keluarganya*”) (gloss: Lim's success makes his family happy). Malay

adjectives can be easily identified if the words are preceded by intensifiers such as “*amat*”, “*sungguh*”, “*sekali*”, “*paling*” and “*agak*”. Yet, it is the opposite in the case of non-adjectives; whereby, over 11% of the words in the hand-tagged Malay corpus are ambiguous [2]. Correspondingly, in recent years, there has been a growing interest in developing data-driven disambiguation applications.

POS-tags can be seen as a disambiguation task since the mapping between words and the tag-space is usually one-to-many [3]. Two possible sources of information can be used to accurately predict the correct POS-tags of a word – contextual information and lexical information. The former is identified based on the different sequences of tags in a sentence. While some POS sequences are common; others are unlikely or impossible. For example, in Malay, prepositional phrases of direction is likely to be followed a verb, a preposition or a noun. On the other hand, the latter is identified based on the semantic value of word itself. For example, the word “*فُكُول* (*pukul*)” (gloss: hit) can either be a verb or a noun. According to [4], the words needs to be analyzed through particular semantic rules to discover whether the meaning is, (first), to hit something, or (second), a special Malay adverbial used to specify the hour in indicating time. However, by utilizing a specific model of statistical and automated learning methods,

features (sequences) of words can be listed without the needs to devise rules that are overcomplicated.

Different approaches have been proposed for the disambiguation tasks of POS-tags. The differences are based on either their internal model, the number of trainings or the information they need to process [3]. In general, these different techniques can be categorized into three major categories: rule-based, statistical-based and transformation-based approaches.

The classical techniques (Rule-based approach) assign its corresponding POS-tags by employing certain lexicography rules. POS-tags which are designed using this approach consists of two stages of architecture [5]. The first stage involves extracting lexicographical data from the dictionary and assigning all probable POS-tags to every word match. The second stage involves employing handcrafted disambiguation rules in order to discover the most appropriate tag for each word.

In the case of automated tagging based on statistical information (statistical-based approach), a lot of different models have been developed and employed. POS tagging based on Hidden Markov Model (HMM) [2], [6]–[8], Maximum Entropy (ME) [9]–[12], Recurrent Neural Network [3], Conditional Random Field [13] and Support Vector Machine [14] are among the models under this category. They are designed based on the statistical occurrences of tag  $n$ -grams and word-tag frequencies which provide the information needed to identify the most probable tag sequence [3].

Transformation-based approaches combine both rule-based and statistical-based approaches. POS-tags based on transformation-based approaches [15] are designed to automatically derive the possible rules directly from the corpora.

Based on two previous studies [2], [11], the performances of POS tagging using HMM and ME models have been compared. HMM for Malay Roman script yielded 67.9% accuracy based on the morphological data gathered; and 94% with TnT. On the other hand, ME with *MaxEnt* and SVM with *SVMTools* reached the overall accuracies of 96% and 99.23% respectively [11]. In another similar study on Bahasa Indonesia [16], an investigation on ME and CRF have been done. ME gives better results (in terms of accuracy) in comparison to CRF. ME recorded an accuracy of 97.57% while CRF recorded an accuracy of 91.15% for two tag sets containing 37 and 25 POS-tags [16].

The objective of this study is to investigate and identify the most appropriate approach for the disambiguation tasks in Malay Jawi POS-tags. Our study is based on the specific contextual information which are related in the Jawi script of Malay corpora.

In Section 2, the standard probabilistic model for POS-tags is presented and discussed. In section 3, ME-based probabilistic model for Malay Jawi will be presented. In section 4, The NUWT Corpus which is used for training and testing POS tagger [17] is

discussed. In section 5, the related contextual information for words and its neighbouring words in the Malay Jawi script will be discussed morphologically (in terms of their suffix/prefix features). In section 6, the training procedure and parameter-setting of ME-based probabilistic model is explained thoroughly. In section 7, the end results and comparative analysis with other methods are presented; and in Section 8, discussions and the conclusion on the findings will be remarked.

## 2.0 POS-TAGS PROBABILISTIC MODEL

A probabilistic model employs POS-tags through the conditional probabilities given by the surrounding contextual features; whereby these probability values are obtained from a manually-tagged corpus [3]. Let  $T = \{t_1, t_2, \dots, t_{|T|}\}$  be a set of POS-tags and  $\Omega = \{w_1, w_2, \dots, w_{|\Omega|}\}$  be the vocabulary of the application. The goal is to find the sequence of POS-tags that maximizes the probability associated to a sentence  $w_1^n = w_1 w_2 \dots w_n$ , i.e.:

$$\hat{t}_1^n = \underset{t_1^n}{\operatorname{argmax}} P(t_1^n | w_1^n). \quad (1)$$

Using Bayes' theorem, the problem is reduced to:

$$\hat{t}_1^n = \underset{t_1^n}{\operatorname{argmax}} P(w_1^n | t_1^n) P(t_1^n). \quad (2)$$

Estimating the values of these parameters can be time consuming since some levels of assumptions are needed – in order to simplify the computational process of the expression [3], [18]. For these models, it is assumed that words are independent of each other; and a word's identity only depends on its tag. Correspondingly, we would be able to obtain this lexical probability,

$$P(w_1^n | t_1^n) \approx \prod_{i=1}^n P(w_i | t_i). \quad (3)$$

Another assumption establishes that the probability of one tag to appear only depends on its predecessor tags,

$$P(t_i^n) \approx \prod_{i=1}^n P(t_i | t_{i-1}, t_{i-2}, \dots, t_{i-k+1}), \quad (4)$$

if a  $k$ -gram class is able to obtain the contextual probabilities.

With these assumptions, a typical probabilistic model following equations (2), (3) and (4) is expressed as follows:

$$\hat{t}_1^n \approx \underset{t_1^n}{\operatorname{argmax}} \prod_{i=1}^n P(w_i | t_i) P(t_{i-1}, t_{i-2}, \dots, t_{i-k+1}). \quad (5)$$

whereby  $\hat{t}_1^n$  is the best estimation of POS-tags for the given sentence  $w_1^n$ . Nonetheless, two limitations on the probabilistic model are identifiable: (1) it does not model long-distance lexical relationships, (2) the

contextual information takes into account the context on the left while the context on the right is not considered [3].

### 3.0 MAXIMUM ENTROPY MODEL FOR MALAY JAWI POS-TAGS

Maximum Entropy (ME) belongs to the family of classifiers known as the exponential or log-linear classifiers [5]. ME is designed to work by extracting some set of features from the input, combining them linearly (multiplying each by a particular weight and then add them up), and then, using this sum as an exponent [5]. This method allows high flexibility in utilizing contextual information and assigns an appropriate tag based on a probability distribution. The probability distribution should have the highest entropy values found on the training corpus and it must be in accordance to certain conditional values. Correspondingly, ME models the POS-tags task as:

$$p(h, t) = \pi \mu \prod_{j=1}^k \alpha_j^{f_j(h, t)} \quad (6)$$

where  $h$  is a 'history' of observation and tag sequences,  $t$  is a tag,  $\mu$  is a normalization constant,  $f_j(h, t)$  is the feature functions with binary values of 0 and 1, and  $\mu$  and  $\alpha_1, \dots, \alpha_k$  are model parameters [16].

The model parameters must be set in a specific value in order to maximize the entropy of the probability distribution; and additionally, the entropy is subjected to the constraints imposed by the value of the  $f_j$  feature functions from the training data [16]. The Generalized Iterative Scaling (GIS) algorithm, Improved Iterative Scaling (IIS) and the optimized version Megam commonly trained these parameters. According to [16], [19], the underlying philosophy is to choose the model that makes the fewest assumptions about the data whilst still remaining consistent with it.

### 4.0 NUWT CORPUS

The NUWT corpora sources were gathered from three different genres of documents – standard written and conversational Malay, Malay narratives and Malay translation of *Quranic* verses. The NUWT corpora are written using Jawi-specific-Buckwalter code [20]. The first source is an annotated corpus named the "Malay corpus". It contains 21 tags and 18,135 tokens with 1,381 words that have ambiguous tags. The corpus was originally prepared by [2] using the Dewan Bahasa dan Pustaka (DBP) tag set; and was written using the Roman script. The second source is a grammatical corpus named the "Malay corpus UKM-DBP"; and is a collection of story books with 12,304 words. The corpus was developed [21] according to the DBP tag set and was also originally written in the Roman script. It has five main tags – with respective elaborated sections for each main tag [17]. The third source is from the "Quranic Malay written in Jawi character Corpus" [22] which is an unannotated text

of Quranic translations. It contains a collection of 114 chapters with 157,388 words. The corpus is written in Jawi standard Unicode (UTF8).

### 5.0 CONTEXTUAL INFORMATION

The training corpus was partitioned into ten parts of equal size. Fundamentally, the words that appeared in each partition played the functions of a testing corpus - enabling the dimensions of the feature sets to be reduced. Additionally, the technique of 10-fold cross-validation was used with 9 different models. From the cross-validation technique, the contextual information can thus, be extracted and concluded with the calculation of an average accuracy. Table 1 shows number of tokens for each fold.

**Table 1** Number of tokens in training corpus for 9 models

Model	Number of tokens	
	Malay corpus	Malay corpus UKM-DBP
1	16,322	11,815
2	14,508	10,502
3	12,695	9,188
4	10,881	7,877
5	9,068	6,564
6	7,254	5,251
7	5,441	3,939
8	3,627	2,626
9	1,814	1,313

Based on the previous work [2] and Jawi rules [23], [24], we consider several types of features, which is likely suitable for Malay Jawi.

### 5.1 AFFIX FEATURES

Affix features are the simplest type of features in Malay language. According to [22], for Malay language, a derived word can be described as a combination of a prefix, a circumfix, a suffix or an infix with a root word. Table 2 exemplifies the differences between the spelling rules for the suffix "+an" in the Roman and the Jawi scripts respectively.

**Table 2** The Roman and Jawi spelling rules for suffixes

Jawi	Roman Scripts
أن+	+an
ن+	+an
ءن+	+an
ان+	+an

Source: [22]

These features are likely to be most useful in languages that utilize morphological rules to modify word structures and meanings such as the Malay language. Additionally, the features have been automatically constructed from the training corpus by

recording all prefixes and suffixes up to a certain length. Table 3 shows the affixation rules applied in the context of Jawi script. From the table, the valid length of affixation in the Jawi script for Malay language is up to 4 morphemes on either side of the stem.

**Table 3** Derivative Jawi writing for prefixes and suffixes

Jawi	Roman Scripts	Jawi	Roman Scripts	Jawi	Roman Scripts
+انتيا	<i>anti+</i>	+م	<i>me+</i>	+فَر	<i>per+</i>
+ءوتو	<i>auto+</i>	+مم	<i>mem+</i>	+فولي	<i>poli+</i>
+ب	<i>be+</i>	+من	<i>men+</i>	+فرا	<i>pra+</i>
+بل	<i>bel+</i>	+مغ	<i>meng+</i>	+فرو	<i>pro+</i>
+بر	<i>ber+</i>	+مغ	<i>menge+</i>	+س	<i>se+</i>
+بي	<i>bi+</i>	+ممفر	<i>memper+</i>	+سوب	<i>sub+</i>
+د	<i>di+</i>	+قنچا	<i>panca+</i>	+سوفرا	<i>supra+</i>
+دفر	<i>diper+</i>	+ف	<i>pe+</i>	+سوا	<i>swa+</i>
+دوي	<i>dwi+</i>	+فل	<i>pel+</i>	+تاتا	<i>tata+</i>
+اىكا	<i>eka+</i>	+قم	<i>pem+</i>	+ت	<i>te+</i>
+جورو	<i>ju ru+</i>	+فن	<i>pen+</i>	+تر	<i>ter+</i>
+ك	<i>ke+</i>	+فغ	<i>peng+</i>	+تري	<i>tri+</i>
+مها	<i>maha+</i>	+فغ	<i>penge+</i>	+تونا	<i>tuna+</i>
+ه	<i>+ah</i>	+يسمى	<i>+isme</i>	+ون	<i>+wan</i>
+ات	<i>+at</i>	+كن	<i>+kan</i>	+واتي	<i>+wati</i>
+يه	<i>+iah</i>	+من	<i>+man</i>	+وي	<i>+wi</i>
+ين	<i>+in</i>	+نيتا	<i>+nita</i>		

Source: [23]

## 5.2 NEIGHBOURHOOD FEATURES

In addition to using the current words, the tags of surrounding words can also be used as features [10]. In this section, we contrastingly used the tags of surrounding words as features. A common example from the Malay linguistic rule is that the word following a cardinal number is often a noun or a verb; but infrequently, it can also be followed by an adjective or preposition. We expect these features to be beneficial to the process of classification in languages that heavily use modifiers and word positioning.

## 6.0 SETTING THE ME-BASED POS TAGGERS

Learning ME model can be done via a generalization of the logistic regression learning algorithms. We want to find the parameter,  $w$ , which maximizes the likelihood of  $M$  training samples:

$$\hat{w} = \underset{w}{\operatorname{argmax}} \prod_i^M P(y^{(i)} | x^{(i)}) \quad (7)$$

Experiments on ME-based model is more complex than any other models. Five experiments on features have been conducted using NLTK module - providing the default parameters with the aptitude to select the best features that can maximize accuracy. Three forms of basic parameters used in the NLTK module are algorithms, log likelihood delta, and training iteration. Table 4 shows the default parameters in NLTK.

**Table 4** NLTK default parameters

Parameters	Type
Algorithm	lls
Log likelihood delta	Stop when the repetition log likelihood fixed less than the previous log likelihood
Maximum iteration	100

The set of features developed are shown in Table 5 and 6. The experiments listed are from a series of preliminary experiments which aims to determine the number of adjustments needed to provide the highest accuracy.

**Table 5** Feature setting for the five experiments

Feature Name	Experiment				
	1	2	3	4	5
suf-1-x	√				
suf-2-xx	√	√	√	√	√
suf-3-xxx	√	√	√	√	√
pre-2-xx			√	√	√
pre-3-xxx		√	√	√	√
word-w <sub>a</sub> -1	√	√	√		
tag-T <sub>a</sub> -1				√	√
tag-T <sub>a</sub> -2					√

**Table 6** Description of the feature set

Feature Name	Condition	Meaning
suf-1-x	= word[-1:] == 'x'	Word ends with "x"
suf-2-xx	= word[-2:] == 'xx'	Word ends with "xx"
suf-3-xxx	= word[-3:] == 'xxx'	Word ends with "xxx"
pre-2-xx	= word[0:1] == 'xx'	Word begins with "xx"
pre-3-xxx	= word[0:2] == 'xxx'	Word begins with "xxx"
word-w <sub>a</sub> -1	= word[-1] == w	Previous word is w
tag-T <sub>a</sub> -1	= tag [-1] == T	Previous word has tag T
tag-T <sub>a</sub> -2	= tag [-2] == T	Two word back has tag T

Source: [10]

In subsequent experiments, only prefixes and suffixes with 2 and 3 morphemes are taken into considerations. The average accuracy was found to have increased with the features as displayed in the results section of these experiments.

## 7.0 ME-BASED TAGGER PERFORMANCE

In determining the best features of the corpus, the set features were run to identify the features with the highest average accuracy. NLTK default parameters were used in producing these results. Correspondingly, five experimental results using five sets of features are presented in Tables 7 and 8 for Malay corpus and

Malay corpus UKM-DBP respectively. By using different feature settings, it displays that the third experiment gave higher level of average accuracy on the Malay corpus. Meanwhile, feature setting for the fourth experiment gives higher average accuracy on Malay corpus UKM-DBP.

For comparative and validation purposes, we tested our corpora using the standard HMM probabilistic model. Table 9 shows the contrasts between these two models. Ultimately, this study has identified that ME model provides a higher average accuracy compared to HMM model for both Malay Jawi corpora. It is possible that these results are due to its source of information to accurately predict the correct POS-tags of a word.

A graph for the Malay corpus and the Malay corpus UKM-DBP is then plotted in Figure 1 to compare the values from the two corpora. It is highly probable that the lower level of accuracy in the Malay Corpus UKM-DBP is due to its genre – narratives.

3	70:30	91.38	89.48	66.03	62.49
4	60:40	91.90	90.30	64.25	60.34
5	50:50	91.01	88.21	65.27	62.36
6	40:60	90.06	87.80	62.15	56.71
7	30:70	88.32	87.22	62.04	56.19
8	20:80	84.30	76.65	56.67	54.23
9	10:90	79.94	68.21	56.05	49.33
	Average	<b>89.30</b>	<b>85.60</b>	<b>62.69</b>	<b>58.53</b>

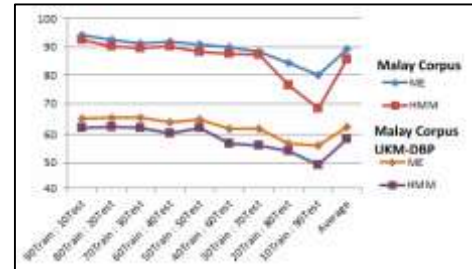


Figure 1 Graph between HMM and ME for Jawi Tagger

Table 7 Experimental results for Malay Corpus

Model (Epoch)	Fold	Malay Corpus				
		1*	2*	3*	4*	5*
1	90:10	88.20	93.83	94.27	94.43	92.17
2	80:20	85.86	92.50	92.48	92.42	92.14
3	70:30	83.22	90.24	91.38	90.87	90.81
4	60:40	82.81	90.57	91.90	91.80	91.90
5	50:50	82.11	89.93	91.01	90.75	90.83
6	40:60	80.81	88.92	90.06	89.92	90.11
7	30:70	79.56	87.01	88.32	88.06	88.22
8	20:80	76.23	82.92	84.30	84.27	84.35
9	10:90	70.94	77.69	79.94	80.61	81.20
	Average	<b>81.08</b>	<b>88.18</b>	<b>89.30</b>	<b>89.24</b>	<b>89.08</b>

Table 8 Experimental results features for Malay Corpus UKM-DBP

Model (Epoch)	Fold	Malay Corpus UKM-DBP				
		1*	2*	3*	4*	5*
1	90:10	60.26	66.12	65.79	65.62	65.38
2	80:20	59.85	64.80	65.99	66.08	65.42
3	70:30	61.28	65.16	66.03	66.03	65.79
4	60:40	59.20	63.66	64.28	64.25	63.16
5	50:50	59.14	64.05	64.96	65.27	62.76
6	40:60	56.12	60.70	61.89	62.15	60.09
7	30:70	55.81	60.60	61.48	62.04	62.48
8	20:80	52.25	55.74	57.08	56.67	57.47
9	10:90	52.55	54.65	56.06	56.05	56.45
	Average	<b>57.39</b>	<b>61.72</b>	<b>62.62</b>	<b>62.69</b>	<b>62.11</b>

\* based on experiment setting

Table 9 Results of the ME highest accuracy and comparison to HMM

Model (Epoch)	Fold	Malay Corpus		Malay Corpus UKM-DBP	
		ME	HMM	ME	HMM
1	90:10	94.27	92.41	65.62	62.45
2	80:20	92.48	90.15	66.08	62.71

## 8.0 DISCUSSION

This study focuses on evaluating the ME model for the development of POS Tags in NLP applications – primarily focusing on its application in the Jawi script of Malay language. The results show that the ME-based model is suitable to be applied to the Malay Jawi script due to its good analytical features on contextual information. Results have also shown that the ME-based model yielded higher accuracy level in comparison to the HMM probabilistic model. The lower level of accuracy in the Malay Corpus UKM-DBP is most probably due to the genre of the corpus.

Based on these findings, a probabilistic model (ME) that can categorize the Jawi-written Malay words into its accurate POS has been identified. For future research endeavours, other Jawi corpora such as the third corpus of NUWT Corpus shall be analyzed for greater reliability and validity. Correspondingly, other derivational words formed through other types of Malay affixations such as circumfix and infix can be added to be part of our future study in NLP applications on the Jawi script of Malay language. Production of the Jawi tagger using ME-based approach will be able to help the intermediate process on NLP onwards.

## Acknowledgement

This material is based upon work supported by the UKM under Grant No. ERGS/1/2013/ICT01/UKM/03/5.

## References

- [1] Biemann, C. 2010. Unsupervised Part-of-Speech Tagging in the Large. *Res. Lang. Comput.* 7(2-4): 101-135.

- [2] Hassan, M., Nazlia, O., and Mohd Juzaidin, A. A. 2011. Statistical Malay Part-of-Speech (POS) Tagger using Hidden Markov Approach. *Proceedings of the International Conference on Semantic Technology and Information Retrieval*. 231–236.
- [3] Zamora-Martinez, F., Castro-Bleda, M. J., Espana-Boquera, S., and Tortajada-Velert, S. 2009. Adding Morphological Information to a Connectionist Part-Of-Speech Tagger. *Current Topics in Artificial Intelligence*, Seville: Springer Berlin Heidelberg. 191–200.
- [4] Knowles, G., and Zuraidah, M. D. 2003. Tagging a corpus of Malay texts, and coping with 'syntactic drift'. *Proceedings of the Corpus Linguistics*. 422–428.
- [5] Jurafsky, D., and Martin, J. H. 2009. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. 2nd Edition. New Jersey, USA: Pearson Education, Inc.
- [6] Hasan, F. M., UzZaman, N., and Khan, M. 2007. Comparison of Different POS Tagging Techniques (N-gram, HMM and Brill's tagger) for Bangla. *Proceeding of the Advances and Innovations in Systems, Computing Sciences and Software Engineering*. 121–126.
- [7] Wicaksono, A. F., and Purwarianti, A. 2010. HMM Based Part-of-speech Tagger for Bahasa Indonesia. *The 4th International MALINDO (Malay and Indonesian Language) Workshop*. 1–7.
- [8] Bar-Haim, R., Sima'An, K., and Winter, Y. 2008. Part-of-speech tagging of Modern Hebrew text. *Nat. Lang. Eng.* 14(02): 223–251.
- [9] Ratnaparkhi, A. 1999. Learning to Parse Natural Language with Maximum Entropy Models. *Mach. Learn.* 34(1–3): 151–175.
- [10] Malecha, G., and Smith, I. 2010. Maximum Entropy Part-of-Speech Tagging in NLTK. 1–10 (unpublished course-related report).
- [11] Hassan, M., Nazlia, O., and Mohd Juzaidin, A. A. 2015. Malay Part Of Speech Tagger: A Comparative Study On Tagging Tools. *Asia-Pacific J. Inf. Technol. Multimed.* 4(1): 11–23.
- [12] Huang, H., and Zhang, X. 2009. Part-of-speech tagger based on maximum entropy model. *Proceeding of the 2nd IEEE Int. Conf. Comput. Sci. Inf. Technol.* 26–29.
- [13] Awasthi, P., Rao, D., and Ravindran, B. 2006. Part Of Speech Tagging and Chunking with HMM and CRF. *Proceedings of NLPAL contest workshop during Nwai '06*. 1–4.
- [14] Sogaard, A. 2010. Simple semi-supervised training of part-of-speech taggers. *Proceedings of the ACL 2010 Conference Short Papers*. 205–208.
- [15] Brill, E. 1995. Transformation-Based Error-Driven Learning and Natural Language Processing: A Case Study in Part-of-Speech Tagging. *Comput. Linguist.* 21(4): 543–565.
- [16] Pisceldo, F., Adriani, M., and Manurung, R. 2009. Probabilistic Part Of Speech Tagging for Bahasa Indonesia. *Proceedings of the Third International MALINDO Workshop, colocated event ACL-IJCNLP*. 1–6.
- [17] Juhaida, A. B., Khairuddin, O., Mohammad Faidzul, N., and Mohd Zamri, M. 2016. NUWT: Jawi-specific Buckwalter Corpus for Malay Word Tokenization. *J. Commun. Inf. Technol* 15:1–25.
- [18] Merialdo, B. Tagging English text with a probabilistic model. 1994. *Comput. Linguist.* 20(2): 155–172.
- [19] Ratnaparkhi, A. 1996. A Maximum Entropy Model for Part-Of-Speech Tagging. *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. 133–142.
- [20] Juhaida, A. B., Khairuddin, O., Mohammad Faidzul, N., Mohd Zamri, M., and Che Wan Shamsul Bahri, C. W. A. C. W. A. 2013. Implementation of Buckwalter transliteration to Malay corpora. *Proceedings of the 13th International Conference on Intelligent Systems Design and Applications*. 213–218.
- [21] Nurul Huda, M. S., Juhaida, A. B., Rafidah, A. K., Nurbaiti, T., and Khalijah, M. N. 2012. Pembangunan korpus cerpen bertag Bahasa Melayu: Analisis Linguistik Korpora. Paper work on *Research, Invention, Innovation & Design (RIID 2012)*. 1–5.
- [22] Suliana, S., Khairuddin, O., Nazlia, O., Mohd Zamri, M., and Hamdan, A. R. 2011. A Malay Stemmers for Jawi Characters. *AI 2011: Advances in Artificial Intelligence*. D. Wang and M. Reynolds, Eds. Perth, Australia: Springer Berlin / Heidelberg. 668–676.
- [23] Hamdan, A. R. 1999. *Panduan menulis dan mengeja Jawi*. Kuala Lumpur: Dewan Bahasa dan Pustaka.
- [24] Ismail, D. 1991. *Pedoman Ejaan Jawi yang Disempurnakan*. Kuala Lumpur: Dewan Bahasa dan Pustaka.