# ABSTRACT

Paraphrase identification is a branch of the NLP study that analyzes the possible paraphrases in two different or more lingual data. microblogging site Twitter is a real example, a news that has information can be rewritten with the same information and different concepts, although the two news have different lexical elements or may have different syntactic structures but have the same meaning can be referred to as paraphrase. In recognizing a form of paraphrase can be done by human evaluation, but the paraphrase evaluation by humans requires a large cost and a longer time, this can be a big problem for developers.

Automatic metric is an automated evaluation engine that uses features that can be used as a lingual extraction to produce a score that can be used as a paraphrase measurement of two sentences compared. In this research, three automatic metric algorithms are used they are BLEU, METEOR, Damerau-Levensthein Edit Distance to identifies paraphrase from the same collected Twitter data. In addition, an analysis of the performance of algorithms by comparing the correlation of human judgment between BLEU, METEOR, Damerau-Levensthein Edit Distance.

From the simulation results conducted in this study, obtained the highest accuracy by using METEOR metric with an accuracy of 0.55 and F1 of 0.76. The second highest value is obtained with BLEU metric with an accuracy value of 0.05 and an F1 value of 0.70. The lowest accuracy score was found on the Distance Distance metric with an accuracy of 0.44 and F1 of 0.30.

*Keywords* : Paraphrase Identification, Tweet, BLEU, Meteor, Edit distance.