

# Clinical Named Entity Recognition in Indonesian Clinical Research Paper using Conditional Random Fields

Rofi Tulus Syifa<sup>1</sup>, Nurmaya<sup>1</sup>, Nova Eka Diana<sup>1</sup>, Yurika Sandra<sup>2</sup>, Wahyu Catur Wibowo<sup>3</sup>, Indra Budi<sup>3</sup>

<sup>1</sup>Faculty of Technology Information, Universitas Yarsi, Jakarta Pusat, Indonesia

<sup>2</sup>Faculty of Medicine, Universitas Yarsi, Jakarta Pusat, Indonesia

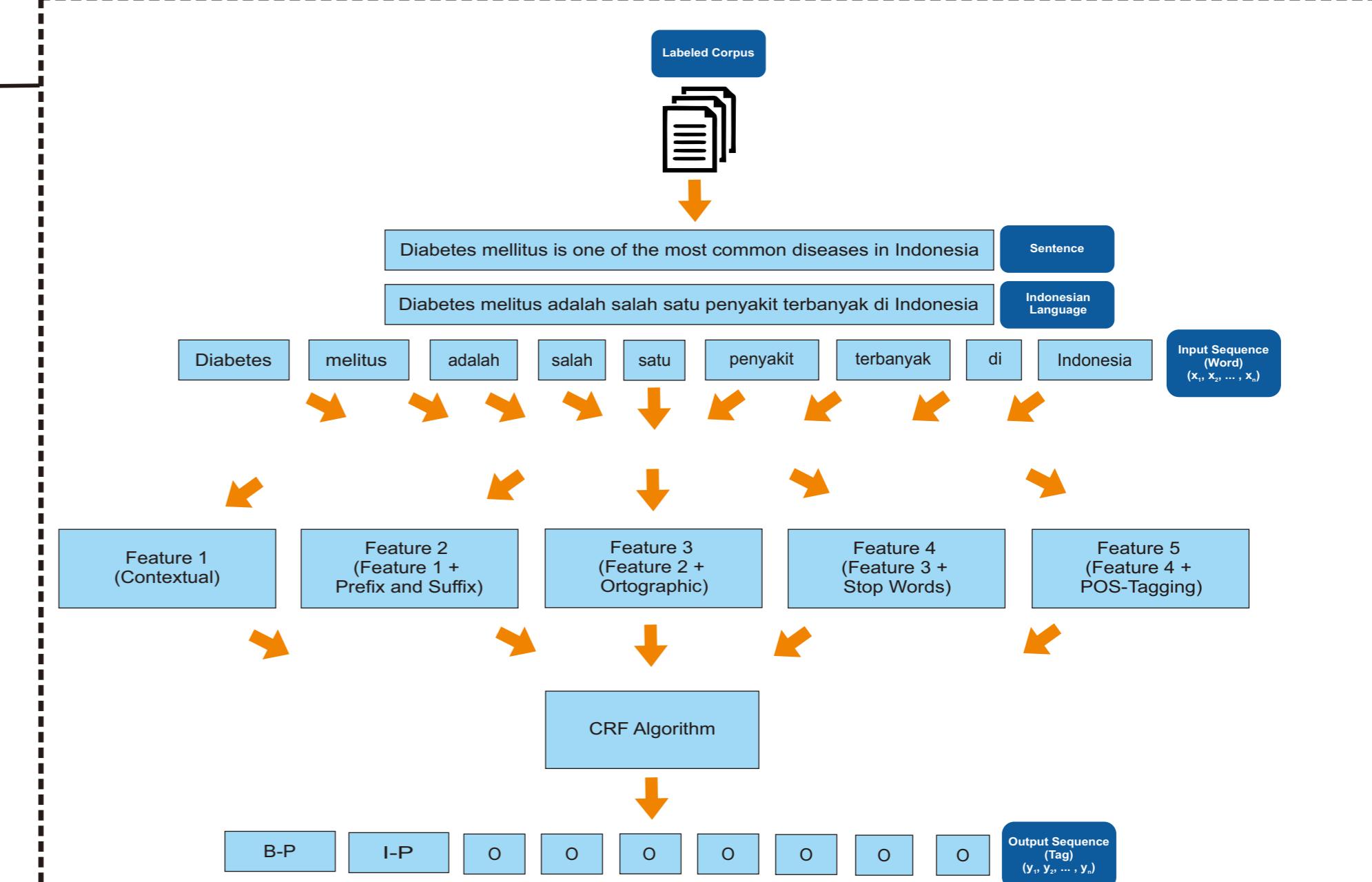
<sup>3</sup>Faculty of Computer Science, Universitas Indonesia, Depok, Indonesia

{rofituluss2@gmail.com, nurmaya@yarsi.ac.id\*}



## Abstract

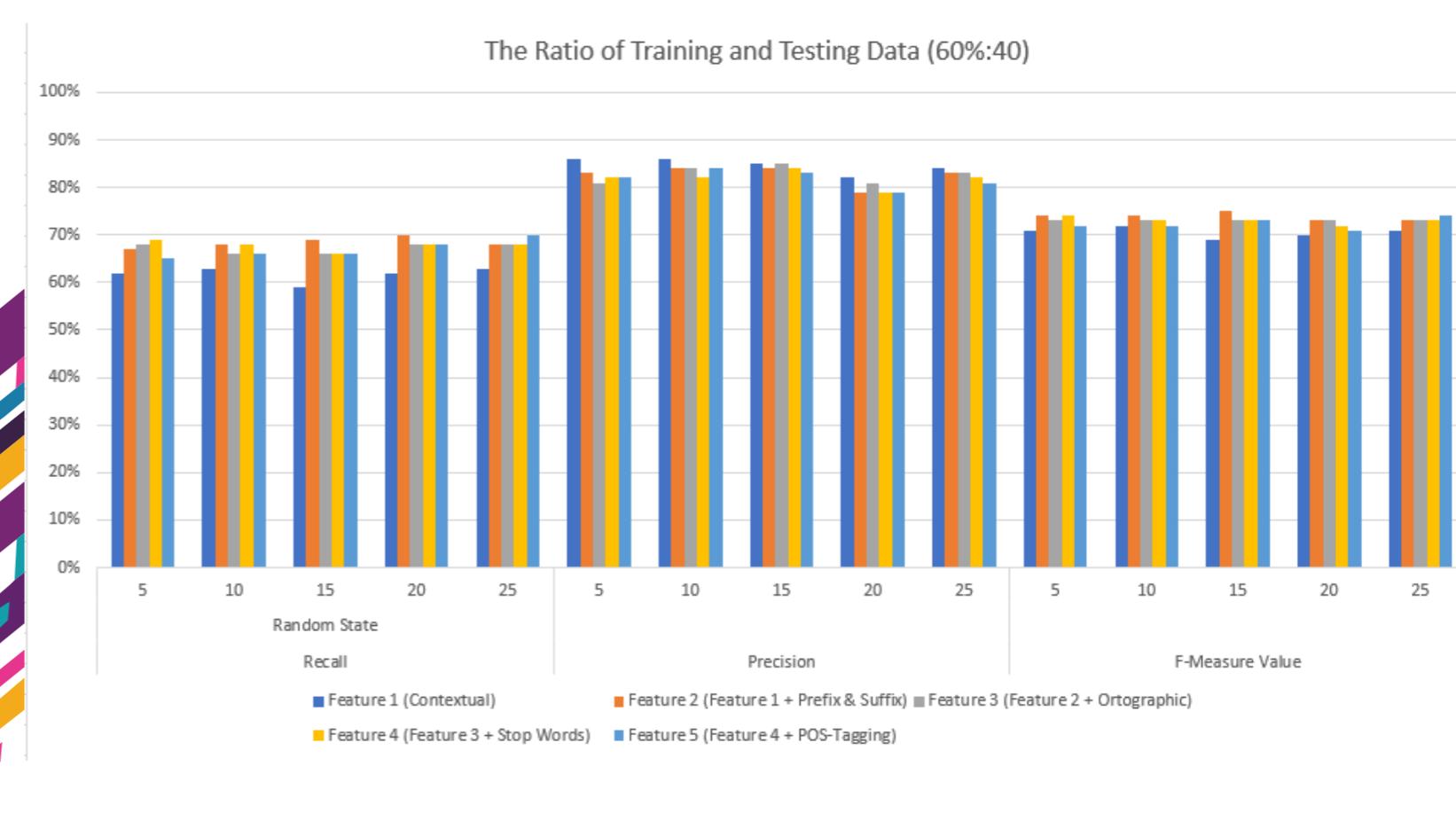
Clinical Research Papers, written in Indonesian Language, contain novelty in clinical knowledge that discusses about patient sign and symptoms, diagnosis, and medical procedures. To assist the readers quickly gain insight into the recent knowledge, it is necessary to extract the entities in Clinical Research Papers. In this study, we have built a Clinical Named Entity Recognition (CNER) model for Indonesian Language using Conditional Random Fields (CRF) with contextual, prefix & suffix, orthographic, stop word, and POS-Tagging features. Using the Indonesian Clinical Corpus annotated in Inside-Outside-Begin (IOB2) format as a data set, we extract seven clinical entities including body part, sign and symptom, disease, physical examination, supporting examination, pharmacological and non-pharmacological treatment. The experiment result shows that our model achieves the best performance in extracting multiple entities with 3 features (contextual, prefix & suffix, and ortographic) with precision of 88%, recall of 73%, and f-measure of 79%.



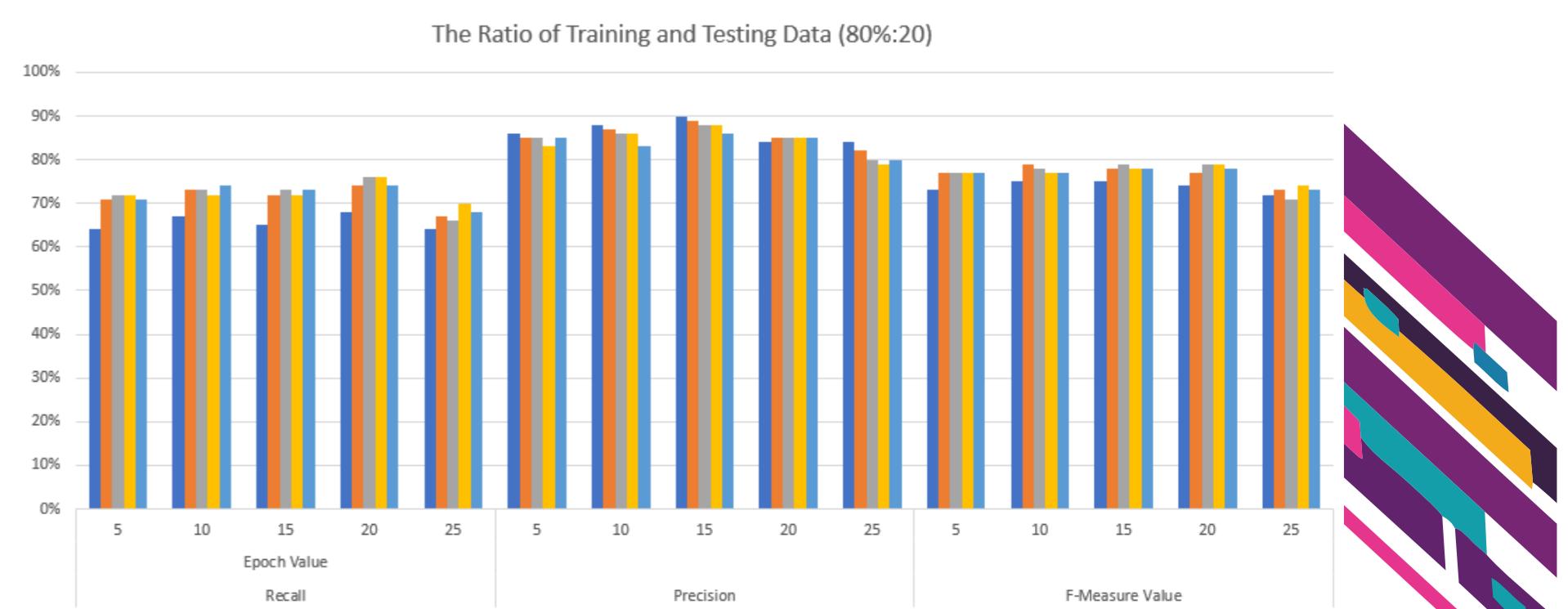
## Result

The experiments conducted in this study are comparing the ratio of training and testing data with random state value and features, comparing the ratio of training and testing data with epoch value. We choose the best CNER Model based on the recall, precision, and F-Measure Value.

A. The ratio of Training and Testing Data (60%:40%) with features and random states



B. The ratio of Training and Testing Data (80%:20) with features and random states



C. The ratio of Training and Testing Data with epoch values

Ratio	Features	Epoch	Recall	Precision	F-Measure Value
60%-40%	Feature 1 (Contextual)	50	61%	84%	70%
	Feature 1 (Contextual)	100	62%	84%	71%
	Feature 1 (Contextual)	150	62%	85%	71%
	Feature 2 (Feature 1 + Prefix & Suffix)	50	63%	82%	73%
	Feature 2 (Feature 1 + Prefix & Suffix)	100	68%	81%	73%
	Feature 2 (Feature 1 + Prefix & Suffix)	150	68%	82%	73%
	Feature 3 (Feature 2 + Orthographic)	50	67%	83%	73%
	Feature 3 (Feature 2 + Orthographic)	100	67%	81%	73%
	Feature 3 (Feature 2 + Orthographic)	150	67%	81%	73%
	Feature 4 (Feature 3 + Stop Words)	50	68%	81%	73%
80%-20%	Feature 4 (Feature 3 + Stop Words)	100	68%	80%	73%
	Feature 4 (Feature 3 + Stop Words)	150	68%	81%	72%
	Feature 5 (Feature 4 + POS-Tagging)	50	67%	80%	72%
	Feature 5 (Feature 4 + POS-Tagging)	100	67%	80%	72%
	Feature 5 (Feature 4 + POS-Tagging)	150	67%	80%	72%
	Feature 1 (Contextual)	50	69%	86%	73%
	Feature 1 (Contextual)	100	69%	86%	73%
	Feature 1 (Contextual)	150	69%	85%	74%
	Feature 2 (Feature 1 + Prefix & Suffix)	50	71%	85%	76%
	Feature 2 (Feature 1 + Prefix & Suffix)	100	73%	85%	77%



Clinical Name Entity Recognition model for Indonesian Language reaches its best performance using CRF algorithm combining with 3 features (contextual, prefix & suffix, and ortographic), the ratio of training & testing data is 80:20, 15 and 50 for random state value and epoch value. The model has precision of 88%, recall of 73%, F-1 measure of 79%.



THE WORKFLOW OF THE SYSTEM IS GIVEN BELOW

