




Universitas Ahmad Dahlan Yogyakarta 46

CEK_ Herman Yuliansyah

-  CEK TURNITIN 3
-  INSTRUCTOR-CEK JURNAL 4
-  Universitas Ahmad Dahlan Yogyakarta

Document Details

Submission ID

trn:oid:::1:2986477510

Submission Date

Aug 20, 2024, 9:32 AM GMT+7

Download Date

Aug 20, 2024, 9:45 AM GMT+7

File Name

98384-361916-1-PB.pdf

File Size

461.8 KB

11 Pages

5,027 Words

28,003 Characters

17% Overall Similarity

The combined total of all matches, including overlapping sources, for each database.

Filtered from the Report

- Bibliography
- Quoted Text

Exclusions

- 50 Excluded Sources

Match Groups

- 54 Not Cited or Quoted 13%**
Matches with neither in-text citation nor quotation marks
- 17 Missing Quotations 4%**
Matches that are still very similar to source material
- 0 Missing Citation 0%**
Matches that have quotation marks, but no in-text citation
- 0 Cited and Quoted 0%**
Matches with in-text citation present, but no quotation marks

Top Sources

- 12% Internet sources
- 11% Publications
- 4% Submitted works (Student Papers)

Integrity Flags

0 Integrity Flags for Review

No suspicious text manipulations found.

Our system's algorithms look deeply at a document for any inconsistencies that would set it apart from a normal submission. If we notice something strange, we flag it for you to review.

A Flag is not necessarily an indicator of a problem. However, we'd recommend you focus your attention there for further review.

Match Groups

- **54 Not Cited or Quoted 13%**
Matches with neither in-text citation nor quotation marks
- **17 Missing Quotations 4%**
Matches that are still very similar to source material
- **0 Missing Citation 0%**
Matches that have quotation marks, but no in-text citation
- **0 Cited and Quoted 0%**
Matches with in-text citation present, but no quotation marks

Top Sources

- 12% Internet sources
- 11% Publications
- 4% Submitted works (Student Papers)

Top Sources

The sources with the highest number of matches within the submission. Overlapping sources will not be displayed.

1	Internet	download.garuda.ristekdikti.go.id	2%
2	Publication	Komang Ayu Triana Indah, I Ketut Gede Darma Putra, Made Sudarma, Rukmi Sari...	1%
3	Internet	e-journals.unmul.ac.id	1%
4	Publication	Chan Sin-wai. "The Routledge Encyclopedia of Translation Technology", Routledg...	1%
5	Internet	ojs.unpkediri.ac.id	1%
6	Student papers	Universitas Mataram	1%
7	Internet	es.scribd.com	1%
8	Publication	Chang Liu, Dongbo Wang, Zhixiao Zhao, Die Hu, Mengcheng Wu, Hai Zhang, Litao...	1%
9	Internet	ejurnal.seminar-id.com	0%
10	Internet	hdl.handle.net	0%

11	Internet	doaj.org	0%
12	Publication	Endah Retnowati, Anik Ghufron, Marzuki, Kasiyan, Adi Cilik Pierawan, Ashadi. "Ch...	0%
13	Publication	Irfan Henuarianto, Ilman Zuhri Yadi, Yesi Novaria Kunang, Susan Dian Purnamas...	0%
14	Publication	Fendi Irfan Amorokhman, Ade Romadhony, Aditya Firman Ihsan. "Indonesian-Kai...	0%
15	Student papers	Universitas Ahmad Dahlan	0%
16	Publication	Zaenal Abidin, Permata Permata, Farida Ariyani. "Translation of the Lampung Lan...	0%
17	Publication	F Rahutomo, A A Septarina, M Sarosa, A Setiawan, M M Huda. "A review on Indon...	0%
18	Publication	"Neural Information Processing", Springer Science and Business Media LLC, 2023	0%
19	Internet	www.nature.com	0%
20	Internet	repository.usd.ac.id	0%
21	Publication	Fatima Zahra El Idrysy, Soufiane Hourri, Ikram El Miqdadi, Assia Hayati, Yassine N...	0%
22	Publication	Mohamed Seghir Hadj Ameur, Farid Meziane, Ahmed Guessoum. "Arabic Machine...	0%
23	Internet	d-nb.info	0%
24	Internet	dokumen.pub	0%

25	Student papers	IUBH - Internationale Hochschule Bad Honnef-Bonn	0%
26	Publication	Ridwan Ilyas, Masayu Khodra, Rinaldi Munir, Rila Mandala, Dwi Widyantoro. "Gen...	0%
27	Publication	Yasir Abdelgadir Mohamed, Akbar Kannan, Mohamed Bashir, Abdul Hakim Moha...	0%
28	Internet	joiv.org	0%
29	Internet	journal1.uad.ac.id	0%
30	Internet	jurnal.untan.ac.id	0%
31	Internet	www.mdpi.com	0%
32	Internet	www.scilit.net	0%
33	Publication	Muthmainah, Nina Zulida Situmorang, Fatwa Tentama. "GAMBARAN SUBJECTIVE ...	0%
34	Publication	Premanand Ghadekar, Neel Malwatkar, Nikhil Sontakke, Nirvisha Soni. "Compara...	0%
35	Publication	Rahmiati Aulia, Diani Apsari, Sri Maharani Budi Haswati, Hana Faza Surya Rusyda ...	0%
36	Internet	ds.libol.fpt.edu.vn	0%
37	Internet	ejournal.st3telkom.ac.id	0%
38	Internet	ejournal.unib.ac.id	0%

39	Internet	ijci.uoitc.edu.iq	0%
40	Internet	itegam-jetia.org	0%
41	Internet	repository.kulib.kyoto-u.ac.jp	0%
42	Internet	www.ijirset.com	0%
43	Publication	Muhammad Khanif Khafidli, Achmad Choiruddin. "Forecast of Aviation Traffic in I..."	0%
44	Publication	Robby Darwis, Herry Sujaini, Rudy Dwi Nyoto. "Peningkatan Mesin Penerjemah St..."	0%
45	Publication	Safitri Nurhaeni, Ruvita Faurina, Ferzha Putra Utama, Kurnia Anggriani. "Sentime..."	0%
46	Publication	Helna Wardhana, I Made Yadi Dharma, Khairan Marzuki, Ibjan Syarif Hidayatulla...	0%
47	Publication	Kyrie Cettyara Eleison, Sari Uli Inggrid Hutahaeen, Sarah Christine Tampubolon, T...	0%

Machine Translation Indonesian Bengkulu Malay Using Neural Machine Translation-LSTM

Bella Okta Sari Miranda¹, Herman Yuliansyah*², Muhammad Kunta Biddinika³

^{1,3}Master of Informatic Universitas Ahmad Dahlan, Yogyakarta, Indonesia

²Department of Informatics, Universitas Ahmad Dahlan, Yogyakarta, Indonesia

e-mail: ¹mirandabella1110@gmail.com, *²herman.yuliansyah@tif.uad.ac.id,

³muhammad.kunta@mti.uad.ac.id

Abstrak

Mesin penerjemah merupakan pengaplikasian dalam Natural Language Processing (NLP) yang berfokus untuk menerjemahkan antar bahasa. Beberapa penelitian terdahulu menggunakan Statistical Machine Translation (SMT) dengan korpus paralel Bahasa Indonesia dan Bahasa Melayu Bengkulu yang berjumlah 3000 data. Namun, SMT memiliki kinerja yang buruk ketika dihadapkan pada jumlah data yang terbatas dan pasangan bahasa yang jarang. Oleh karena itu, penelitian ini bertujuan untuk membangun model mesin penerjemah dari Bahasa Indonesia ke Bahasa Melayu Bengkulu menggunakan pendekatan NMT dengan menggunakan Long Short-Term Memory (LSTM) dan membuat korpus paralel sebanyak 5261 pasangan data Bahasa Indonesia dan Bahasa Melayu Bengkulu. Penelitian ini dilakukan dalam tiga tahapan, yaitu pengumpulan data, prapemrosesan data, pelatihan dan pemodelan, serta evaluasi. Evaluasi kinerja mesin penerjemah dilakukan dengan menggunakan Bilingual Evaluation Understudy (BLEU). Hasil evaluasi menunjukkan bahwa model ini mencapai rata-rata tertinggi sebesar 0.6016332 pada BLEU-1 dan rata-rata terendah sebesar 0.3680788 pada BLEU-4. Hasil ini menunjukkan pertimbangan perbedaan struktur linguistik alami antara Bahasa Indonesia dan Bahasa Melayu Bengkulu dapat disarankan sebagai solusi terbaik untuk menerjemahkan dari Bahasa Indonesia ke Bahasa Melayu Bengkulu.

Kata kunci— Bahasa Melayu Bengkulu, BLEU, NMT, Korpus Paralel, LSTM

Abstract

The machine translator is an application in Natural Language Processing (NLP) that focuses on translating between languages. Several previous research have used Statistical Machine Translation (SMT) with a parallel corpus of Indonesian and Bengkulu Malay totaling 3000 data points. However, SMT performs poorly when confronted with limited data and infrequent language pairs. Therefore, this study aims to build a machine translation model from Indonesian to Bengkulu Malay using an NMT approach with Long Short-Term Memory (LSTM), and to create a parallel corpus of 5261 data pairs between Indonesian and Bengkulu Malay. The research was conducted in three stages: data collection, data preprocessing, training and modeling, and evaluation. The performance of the machine translator was evaluated using the Bilingual Evaluation Understudy (BLEU). The evaluation results show that this model achieved the highest average score of 0.6016332 on BLEU-1 and the lowest average score of 0.3680788 on BLEU-4. These results indicate that considering the natural linguistic structural differences between Indonesian and Bengkulu Malay can be suggested as the best solution for translating from Indonesian to Bengkulu Malay.

Keywords— Bengkulu Malay Language, BLEU, NMT, Parallel Corpus, LSTM

Received June 1st, 2012; Revised June 25th, 2012; Accepted July 10th, 2012

1. INTRODUCTION

Language is a communication tool used to interact and share information between individuals. Apart from functioning as a medium for conveying opinions and arguments to others, language also plays an important role in the social context [1]. Bengkulu Malay has become one of the languages commonly used by Bengkulu residents in various aspects of daily life [2]. From everyday conversations to social and cultural activities, Bengkulu Malay is important in facilitating communication between individuals, groups, and communities in the area. As a language that permeates everyday life, Bengkulu Malay is also a reflection of local cultural identity and heritage, reflecting the plurality and diversity of Bengkulu society [3]. Therefore, the existence and use of the Bengkulu Malay language not only enriches social and cultural life in Bengkulu but also strengthens relations between community members in the region.

Thus, language is not only a means of communication but also facilitates the exchange of ideas and views in the wider community. The plurality of languages throughout the world is the main cause of difficulties in communication between people. Efforts to maintain the sustainability of regional languages involve teaching the younger generation about these languages and encouraging their use in daily activities. Real steps in preserving regional languages can ensure that cultural heritage remains relevant in the era of globalization. Therefore, a medium is needed to overcome the challenges of diversity, and one solution is through the use of machine translation.

Machine translation is a subfield of computational linguistics that falls under the category of Natural Language Processing (NLP) [4]. Machine translation can be defined as a computerized system for translating from one language to another [5]. Previous research utilized Statistical Machine Translation (SMT) with a parallel corpus of Indonesian and Bengkulu Malay languages totaling 3000 data points [6]. However, SMT performs poorly when faced with limited data and sparse language pairs. Common methods used in machine translation development include SMT, Rule-Based Machine Translation (RBMT), and Neural Machine Translation (NMT) [7]. NMT is a translation approach that utilizes LSTM with encoder and decoder architectures [8].

Translation from Indonesian to Sundanese using RNN is an implementation of technology that utilizes a recurrent neural network model to process and translate text between these two languages [9]. Designing a Neural Network-based Translator from Kawi language to Indonesian language using the Microframework Flask involves developing a system that utilizes neural network technology to process and translate text between these two languages [10]. Sentence translation from the Lampung language to Indonesian using NMT with Attention is a process where a computer system utilizes a neural network-based model to automatically translate text from the Lampung language to Indonesian while considering the context and relationships between words in the sentence [11]. An application that translates the Bangka language to Indonesian using NMT based on a website is a system that utilizes neural network technology to automatically translate text between these two languages. With a web-based platform, users can access this application through a web browser to input text in the Bangka language and receive the translation results in Indonesian directly [12]. Testing the accuracy of NMT from Indonesian to Tiochiu Pontianak using the Bahdanau Attention mechanism aims to measure how well this automatic translation system understands and produces appropriate translations within the context of Tiochiu Pontianak language [13].

Research on SMT highlights weaknesses due to separate components such as translation models, language models, and parsing models. This approach requires careful coordination among these components to achieve accurate translation results. In contrast, NMT employs an end-to-end approach in which the entire model is treated as a single integrated unit. This means that NMT integrates the entire process from source text to target text within a single model, without requiring separate components as in SMT. The end-to-end approach can reduce system

complexity and enhance the model's ability to understand and translate text more effectively, as it allows the model to directly learn the relationship between source and target texts.

This research aims to develop a parallel corpus comprising 5261 data points and to develop a machine translation model for Indonesian to Bengkulu Malay using LSTM architecture. With a significant increase in the size of the parallel corpus, it is expected that the resulting machine translation model will have improved capabilities in understanding and translating texts between Indonesian and Bengkulu Malay with high accuracy. Subsequently, the performance of the proposed model will be evaluated using the BLEU method [14][15].

2. METHODS

The research methodology is divided into four parts: data collection, data preprocessing, modeling and training, and evaluation, as shown in Figure 1.

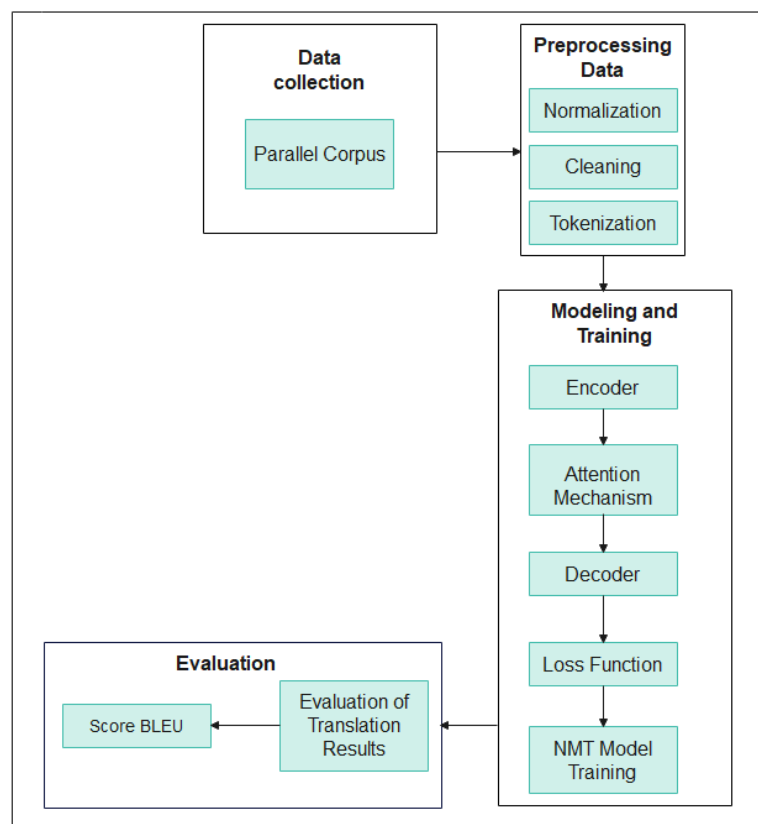


Figure 1. Design of the Indonesian-Bengkulu Machine Translation System

2.1 Data Collection

The data in this study uses a parallel corpus. A parallel corpus is a collection of texts aligned or synchronized in two or more languages [16]. Parallel corpora can be applied to translation processes and expressing ideas in two (or more) different languages, or to compare the characteristics of original (source) texts and translated texts [17]. Parallel corpora are also beneficial for other natural language processing applications such as cross-lingual information retrieval, disambiguation of words, and projection of explanations [6]. The parallel corpus in this study was created independently by manually typing each sentence due to the limited availability of online resources. It comprises 5261 sentences, with Indonesian sentences sourced from online resources and Bengkulu Malay sentences taken from the Bengkulu Malay dictionary available at the local library in Bengkulu City, Bengkulu province, Indonesia. The

research data consists of both languages: Indonesian as the source language and Bengkulu Malay as the target language. The dataset for this research can be seen in Table 1.

Table 1. Parallel Corpus

No	Bahasa Indonesia	Bahasa Melayu Bengkulu
1	Bunga dihalaman rumah dibersihkan semua oleh tukang rumput	Bungo dilaman rumah di babat galo dekek tukang rumput
2	Bapak ini sangat suka dengan sayur perut sapi	Gaek lanang iko suko nian gulai sapi babat
3	Dia menikahkan anak nya sangat besar-besaran	Nyo nikahi anak nyo berarat-arat nian
4	Sakitnya sekarang sudah bertambah parah	Sakitnyo kini lah bertambah arat
5	Sudah sebesar ini saya belum juga menghapal perkalian	Lah segedang iko ambo belum jugo apal perkalian
6	Mengapa dia tidak mau menyingkir dari rumah itu	Ngapo nyo idak endak berasak dari rumah itu
7	Mengapa bau baju ini sangat kurang sedap	Ngapo baun baju ko apak nian
8	Kamu jangan pernah membawa fitnah	Kau janganlah membaok apat
9	Karpet ini saya beli di mekkah	Ambal iko ambo beli di mekkah
10	Saya lupa mengunci pintu rumah	Ambo lupu ngunci pintu rumah
...		
5261	Api pemadam cepat membatasi api agar tidak merambat.	Api pemadam cepek ngebatasi api biar idak merambek

The parallel corpus data of Indonesian and Bengkulu Malay is provided to facilitate researchers or other users access linguistic resources relevant to their research. This aims to support ease of access and utilization of information in linguistic contexts that are pertinent to research needs. The data can be accessed at the following URL: (<https://data.mendeley.com/datasets/tnk42hhjk/1>)

2.2 Preprocessing Data

Data preprocessing is a stage where initially unstructured data is transformed into structured data. This is done by optimizing the data through customized processes. Data preprocessing is performed to prepare the data for use in learning processes. By preprocessing, raw data in the dataset is processed and adjusted so that it can be processed by the algorithms to be used. In this study, preprocessing stages include normalization, data cleaning, and tokenization.

Normalization of data in the preprocessing stage involves removing unnecessary Unicode characters, converting text to lowercase, and inserting spaces before and after specific punctuation marks such as . ? , ! Data cleaning is performed to replace multiple consecutive spaces with a single space, replace non-letter characters and punctuation marks with spaces, and remove any leading or trailing spaces from the string [18]. Tokenization in the preprocessing process separates words from each other or punctuation marks [19].

2.3 Modeling and Training NMT

During the training and modeling stage of NMT, there is a process of mapping from the input sequence of words to the output sequence of words. However, in practice, the translated input word sequences often differ from the output sequences. Moreover, the length of the input sequence is not always the same as the length of the output sequence, and the positions of the generated sequences can also vary. Sequence to sequence (SeqToSeq) is known as a model for

handling mapping from input sequences to output sequences. NMT utilizes the SeqToSeq model framework, which consists of two main networks: the encoder and the decoder [20]. An illustration of how to connect the encoder and decoder models is show in Figure 2.

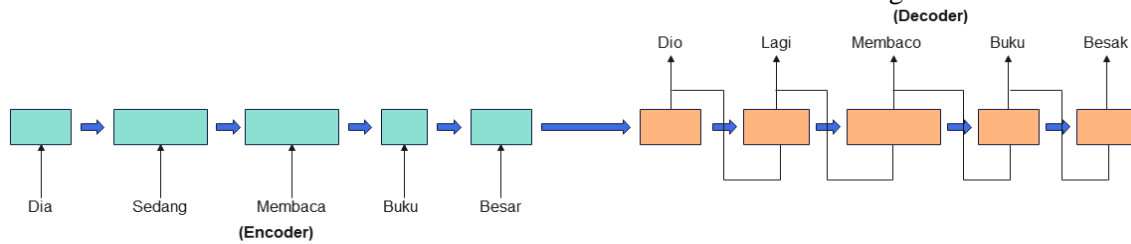


Figure 2. Encoder-Decoder architecture model

Figure 2 illustrates the relationship between the encoder and decoder models in the context of Indonesian and Bengkulu Malay languages, using the example sentence in Indonesian "Dia Sedang Membaca Buku Besar" (He Is Reading Big Book). It depicts the encoder model receiving input in Indonesian and generating a representation that the decoder uses to produce a corresponding translation in Bengkulu Malay.

The first process undertaken is the encoder. The encoder process begins by converting the source language into a numeric representation, where each word is mapped to a unique integer. This is necessary because computers can only process and understand numeric data, not strings or characters. This process involves building a dictionary data structure that maps each word to a unique integer and allows for the reverse conversion from integers back to corresponding words [21].

After the encoder process, the next step involves the connection process between the encoder and decoder, which plays a crucial role in learning and determining scores to align input sequences [22][23]. This research applies the attention mechanism as an essential link between the encoder and decoder in the context of developing a machine translation system. This mechanism is crucial in learning and determining the adjustment of input sequences that are most relevant to the translation context being addressed. The primary focus of the attention mechanism is on calculating alignment scores, which leads to selecting important information from the source text (encoder) to use in generating the output text (decoder). The attention mechanism used in this study is the Luong attention mechanism, which implements three different score calculation methods: dot product, general, and concat. Each method has its advantages and characteristics in capturing the complex relationships between words and phrases in the two translated languages.

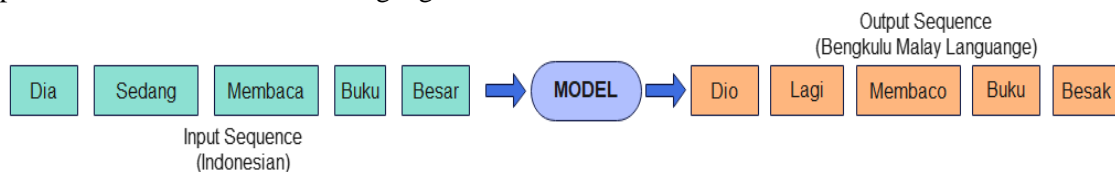


Figure 3. Sequential Model for Indonesian-Bengkulu Malay Machine Translation

The next process is the decoder. The decoder involves steps to generate output, which is predicting words at each time step based on the processes that have been undergone previously [24]. Similar to the encoder, the decoder also consists of two layers, but it utilizes an attention layer to obtain output. The attention layer provides information in the form of an attention vector that stores attention weights from all encoder hidden states, thereby assisting the decoder in focusing at each time step. The output generated at each decoder time step functions as input to the next context vector [25]. The model trains the dataset representing words in the form of unique integers, which are then converted into tensors. Therefore, the sparse categorical cross-entropy module is used because this dataset has many categories (number of source and target words) [26].

The next step is to proceed with training the NMT model. The steps to train the NMT model involve using training data and applying specific techniques over a certain period to enable the model to make predictions. During the training process, input tensors representing input sentences from the encoder are used. This generates encoder outputs containing hidden states from all time steps. Subsequently, the last hidden state of the encoder is taken as the initialization for the initial hidden state of the decoder [27].

Word embedding is performed as a process involving converting words into indices using sequence encoding methods. The analysis results indicate that the embedding layer accepts inputs with a maximum size of 18, representing vectors in the source language, Indonesian. The final layer generates output with a size of 10, representing the maximum vector size for the target language, Malay Bengkulu. Therefore, both languages have the same maximum length in this text-processing context. The model summary can be viewed in Figure 4.

Layer (type)	Output Shape	Param #
input_1 (InputLayer)	[(None, 19)]	0
input_2 (InputLayer)	[(None, 18)]	0
embedding (Embedding)	(None, 19, 256)	834560
embedding_1 (Embedding)	(None, 18, 256)	732160
lstm (LSTM)	[(None, 19, 256), (None, 256), (None, 256)]	525312
lstm_1 (LSTM)	[(None, 18, 256), (None, 256), (None, 256)]	525312
dot (Dot)	(None, 18, 19)	0
activation (Activation)	(None, 18, 19)	0
dot_1 (Dot)	(None, 18, 256)	0

Figure 4. Model Summary

This research has been tailored by implementing an embedding layer to represent sentence lengths in Indonesian as input. At the same time, the vocabulary size in Malay Bengkulu is used as output. The chosen batch size in this research is 64, decided based on the available RAM capacity of the system to ensure efficiency and stability during the model training process. The optimizer used is the Adam optimizer, known for its reliability in optimizing weights in model development contexts. The activation function chosen is Softmax, and a dropout layer with a value of 0.2 has been added to reduce the risk of overfitting. The detailed parameters applied to the model architecture in this study can be found in Table 2.

Table 2. Parameter model architecture

Parameter	Value
Batch Size	64
Loss	Categorical Crossentropy
Dropout	0.2

Activation	Softmax
Optimizer	Adam

The training data used in this study consists of 4209 parallel corpora. This dataset serves as the primary dataset used to train the model or system developed in the research. On the other hand, the testing data consists of 1052 parallel corpora used to evaluate the performance and accuracy of the trained model or system.

2.4 Evaluation

The evaluation stage in translation sentences is concluded by automatic evaluation to measure their quality. Automatic testing is carried out using a consistent approach to increasing the number of epochs. In the evaluation process, after the model is updated with training data, testing is conducted by increasing the number of epochs to observe changes in the accuracy and consistency of translation results. This method allows for identifying how increasing the number of epochs affects the overall performance of the model, both in terms of accuracy and its ability to adapt to new data.

In this research, the automatic testing process is conducted using the BLEU metric to evaluate the translation results from 1-gram to 4-gram. BLEU is an algorithm that evaluates the quality of machine translation from one natural language to another [28]. BLEU measures the modified n-gram precision score between the automatic translation output and the reference translation and utilizes a constant called brevity penalty [29]. The BLEU score is obtained by multiplying the brevity penalty with the geometric mean of the modified precision scores. A high BLEU score is achieved when the length of the translated sentence approaches that of the reference sentence, and the translated sentence matches the reference sentence in terms of words and their order [30]. The BLEU formula is shown in Equation 1.

$$BP_{BLEU} = f(x) = \begin{cases} 1, & \text{if } c > r \\ e^{(1-\frac{r}{c})}, & \text{if } c \leq r \end{cases} \quad (1)$$

$$P_n = \frac{\sum_{C \in \{Candidates\}} \sum_{ngram \in \{countclip(ngram)\}}}{\sum_{C \in \{Candidates\}} \sum_{ngram \in \{countclip(ngram)\}}} \quad (2)$$

$$BLUE = BP \cdot \exp \sum_{n=1}^N W_n \cdot \log P_n \quad (3)$$

Keterangan:

BP = Brevity Penalty

c = Number of words in the automatic translation result (candidate)

r = Number of references (reference)

W_n = 1/N (N for BLUE which is 4)

P_n = Number of n-grams in the translation output that match the reference divided by the number of n-grams in the translation output.

3. RESULTS AND DISCUSSION

Several tests in this study were conducted to assess the effectiveness of various aspects in the development of the NMT model using the encoder-decoder method. Testing involved evaluating the performance of the translation engine by trying various numbers of epochs, ranging from 10 to 100. This research considers not only the final evaluation results but also the time required to train the model. The batch size selected in this study was 64, chosen based on the available RAM capacity of the system to ensure efficiency and stability during the model

training process. The optimizer used was the Adam optimizer, known for its reliability in optimizing weights in the context of model development.

The evaluation is conducted by comparing the BLEU scores generated at the end of the analysis. This comparison aims to assess the performance of the machine translation system in the context of relevant applications. The BLEU metric is used as an indicator to measure how well the machine translation results approximate the reference or expected outcomes.

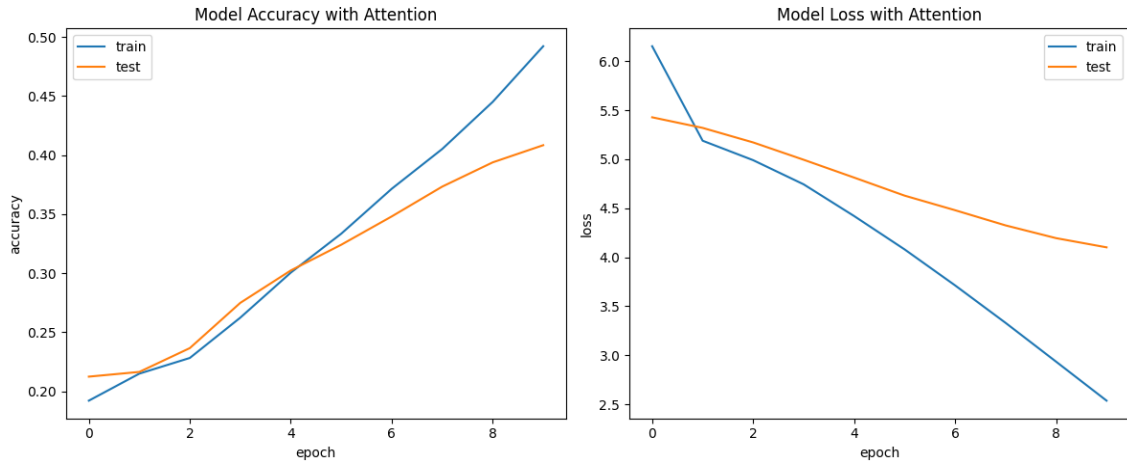


Figure 5. Model with Epoch 10

In Figure 5, evaluation was performed using a limited number of epochs, specifically only 10 epochs. The results of this evaluation showed that the BLEU score obtained was 0.324243. BLEU score is used to measure how well machine translation results align with human reference translations. With this score, it can be estimated how well the model captures and reproduces accurate sentence structures in translation, despite the score being relatively low. This evaluation provides an initial overview of the model's performance under the set epoch limitation.

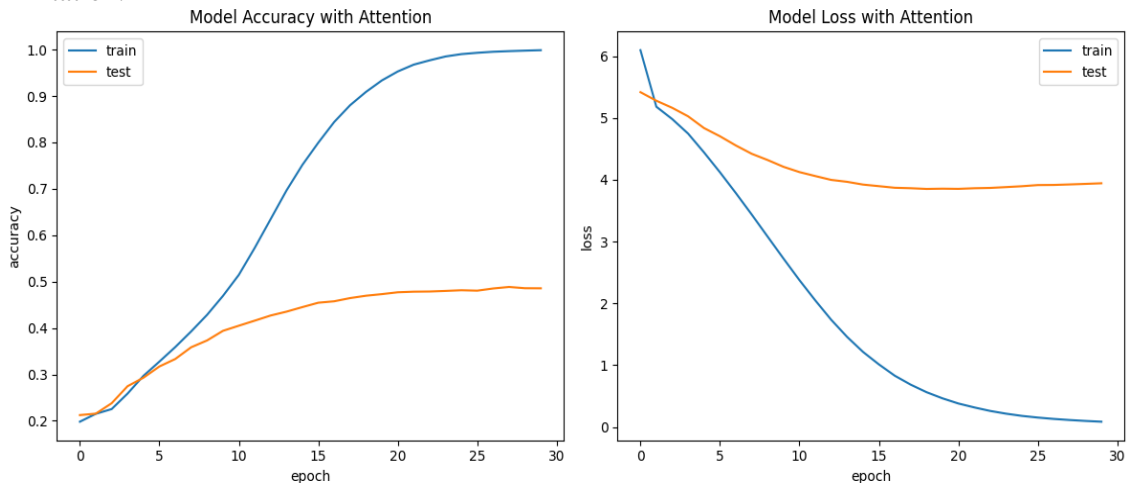


Figure 6. Model with Epoch 30

The Figure 6 shows the results of testing using the Attention mechanism with 30 epochs. The use of the Attention mechanism in this test helps to improve understanding of how the model allocates attention to various parts of the input data, thereby enabling more accurate predictions.

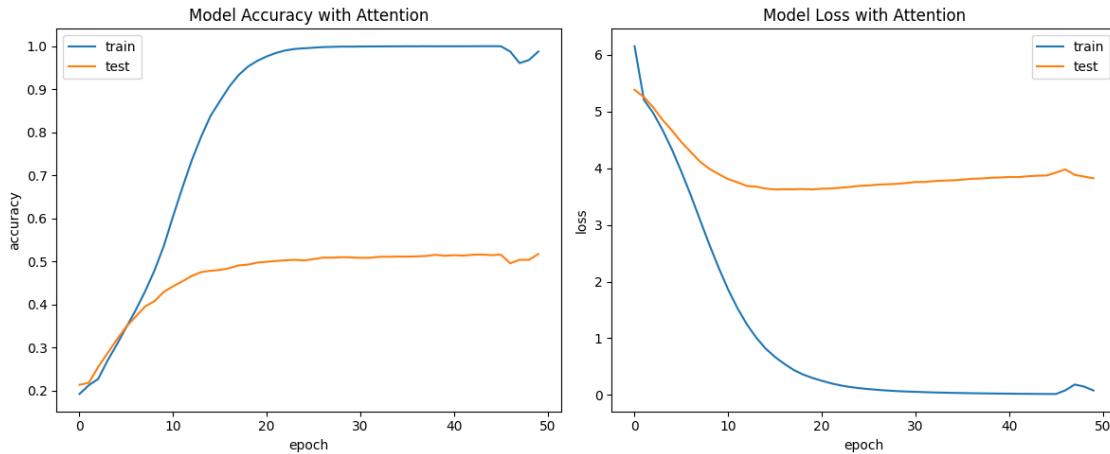


Figure 7. Model with Epoch 50

The Figure 7, shows that the model’s loss reaches its lowest point at epoch 50. After reaching this point, the loss gradually increases. This indicates that the model tends to overfit, where its performance on training data is better than on validation or test data.

Additionally, the observed difference in BLEU scores also indicates this phenomenon. The BLEU score for training data is 0.668649, showing that the model performs relatively well on the data used to train it. However, the lower BLEU score on test data, which is 0.342092, indicates a decrease in performance when applied to unseen data.

Table 3. Evaluating the model with BLEU Score

Metric	Scenario 1 Epoch 10	Scenario 2 Epoch 30	Scenario 3 Epoch 50	Scenario 4 Epoch 70	Scenario 5 Epoch 100	Average
BLEU-1	0.324243	0.647730	0.668649	0.676953	0.690591	0.6016332
BLEU-2	0.179242	0.557949	0.586816	0.595997	0.615129	0.5070266
BLEU-3	0.117397	0.518461	0.549068	0.559067	0.581134	0.4650254
BLEU-4	0.046645	0.413601	0.443655	0.455900	0.480593	0.3680788

This research involved five carefully conducted experiments, each experiment result is meticulously recorded in Table 2 to evaluate the BLEU metric scores. The BLEU scores in the testing table indicate values for 1-grams, 1-2 grams, 1-3 grams, and 1-4 grams respectively. BLEU-1 (which measures unigram similarity between the model translation and the reference) shows the highest average score at 0.6116332, BLEU-2 shows an average of 0.5070266, and BLEU-3 shows an average of 0.4650254. BLEU-4 (which measures similarity up to four-token n-grams) shows the lowest average score at 0.3680788.

The minimal difference between the scores on training and test data indicates that the model did not experience significant overfitting to the training data. The goal of these experiments was to gain a deep understanding of the model’s performance in text translation. Thus, a comprehensive analysis of the resulting data can provide deeper insights into the effectiveness and accuracy of the approach used.

4. CONCLUSIONS

This research develops a machine translator from Indonesian to Bengkulu Malay using an NMT approach that integrates LSTM. The parallel corpus used in this system comprises 5261 pairs of Indonesian and Bengkulu Malay data. An LSTM encoder-decoder model is utilized to evaluate BLEU scores, achieving the highest average of 0.6016332 for BLEU-1 and 0.3680788 for BLEU-4. Considering the natural linguistic structural differences between Indonesian and Bengkulu Malay, this LSTM is recommended as the best solution for translating

Title of manuscript is short and clear, implies research results (First Author)

9
8

27
26

1

from Indonesian to Bengkulu Malay. Recommendations for future research include expanding the dataset to enhance prediction variability in NMT models. This will facilitate the recognition of various words during the language model training phase and support the implementation of Transformer Neural Network as a language model in machine translation.

ACKNOWLEDGEMENTS

This research was supported by the Direktorat Riset, Teknologi, dan Pengabdian Masyarakat (DRTPM) Direktorat Jenderal Pendidikan Tinggi, Riset, dan Teknologi, Kementerian Pendidikan, Kebudayaan, Riset dan Teknologi under Penelitian Tesis Magister (Master's Thesis Research) with grant number: 061/PTM/LPPM-UAD/VI/2024 (15 Juni 2024).

REFERENCES

- [1] F. Senovil, "Morfofonemik Bahasa Melayu Bengkulu," *KLITIKA J. Ilm. Pendidik. Bhs. dan Sastra Indones.*, vol. 2, no. 2, pp. 165–178, 2020, doi: <https://doi.org/10.32585/klitika.v2i2.1037>.
- [2] N. H. M. Ningsih, D. E. C. Wardhana, and S. Supadi, "Derivasi Bahasa Melayu Bengkulu," *J. Ilm. KORPUS*, vol. 4, no. 2, pp. 224–230, 2020, doi: 10.33369/jik.v4i2.8361.
- [3] J. Zakaria, I. Yuniati, and E. F. Wijaya, "Implikatur Tegur Sapa Dalam Bahasa Melayu Bengkulu," *Lit. J. Bahasa, Sastra dan Pengajaran*, vol. 1, no. 2, pp. 74–78, 2021, doi: <https://doi.org/10.31539/literatur.v1i2.2401>.
- [4] M. Stasimioti, V. Sosoni, D. Mouratidis, and K. Kermanidis, "Machine Translation Quality: A comparative evaluation of SMT, NMT and tailored-NMT outputs," *Proc. 22nd Annu. Conf. Eur. Assoc. Mach. Transl. EAMT 2020*, pp. 441–450, 2020.
- [5] A. Garg and M. Agarwal, "Machine Translation : A Literature Review."
- [6] D. Soyusiawaty and B. O. S. Miranda, "Statistical Machine Translation from Indonesian to Regional Languages in Indonesia," *Int. J. Comput. Appl.*, vol. 184, no. 49, pp. 18–23, 2023, doi: 10.5120/ijca2023922603.
- [7] F. Rahutomo, A. A. Septarina, M. Sarosa, A. Setiawan, and M. M. Huda, "A review on Indonesian machine translation," *J. Phys. Conf. Ser.*, vol. 1402, no. 7, 2019, doi: 10.1088/1742-6596/1402/7/077040.
- [8] Z. Tan, S. Wang, Z. Yang, G. Chen, and X. Huang, "Neural machine translation : A review of methods , resources , and tools," *AI Open*, vol. 1, no. October 2020, pp. 5–21, 2021, doi: 10.1016/j.aiopen.2020.11.001.
- [9] Y. Fauziyah, R. Ilyas, and F. Kasyidi, "Mesin Penerjemah Bahasa Indonesia-Bahasa Sunda Menggunakan Recurrent Neural Networks," *J. Teknoinfo*, vol. 16, no. 2, pp. 313–322, 2022, doi: <https://doi.org/10.33365/jti.v16i2.1930>.
- [10] I. G. A. Budaya, M. W. A. Kesiman, and I. M. G. Sunarya, "Perancangan Mesin Translasi berbasis Neural dari Bahasa Kawi ke dalam Bahasa Indonesia menggunakan Microframework Flask," *J. Sist. dan Inform.*, pp. 94–103, 2022.
- [11] Z. Abidin, A. Sucipto, and A. Budiman, "Penerjemahan Kalimat Bahasa Lampung-Indonesia Dengan Pendekatan Neural Machine Translation Berbasis Attention Translation of Sentence Lampung-Indonesian Languages With Neural Machine Translation Attention Based," *J. Kelitbangan*, vol. 06, no. 02, pp. 191–206, 2018.
- [12] F. Razsiah, A. Josi, and S. Mubaroh, "Aplikasi Penerjemah Bahasa Bangka Ke Bahasa Indonesia Menggunakan Neural Machine Translation Berbasis Website," *J. Inov. Teknol. Terap.*, vol. 1, no. 1, pp. 68–76, 2023, doi: 10.33504/jitt.v1i1.67.
- [13] L. B. San and H. Sujaini, "Uji Nilai Akurasi pada Neural Machine Translation (NMT) Bahasa Indonesia ke Bahasa Tiochiu Pontianak dengan Mekanisme Attention," vol. 9, no. 3, pp. 362–370, 2023. doi: <https://dx.doi.org/10.26418/jp.v9i3.63346>
- [14] D. A. Sulistyio, A. P. Wibawa, D. D. Prasetya, and F. A. Ahda, "LSTM-Based Machine

- Translation for Madurese-Indonesian,” *J. Appl. Data Sci.*, vol. 4, no. 3, pp. 190–199, 2023, doi: 10.47738/jads.v4i3.113.
- [15] K. Dedes *et al.*, “Neural Machine Translation of Spanish-English Food Recipes Using LSTM,” *Int. J. Informatics Vis.*, vol. 6, no. June, pp. 290–297, 2022, [Online]. Available: www.joiv.org/index.php/joiv. doi: <https://dx.doi.org/10.30630/joiv.6.2.804>
- [16] Q. A. Agigi and A. A. Suryani, “Statistical Machine Translation Muna to Indonesia Language,” *J. Tek. Inform. dan Sist. Inf.*, vol. 8, no. 4, pp. 2173–2186, 2021, doi: 10.35957/jatisi.v8i4.1149.
- [17] M. S. Alam and A. A. Suryani, “Minang and Indonesian Phrase-Based Statistical Machine Translation,” *J. Informatics Telecommun. Eng.*, vol. 5, no. 1, pp. 216–224, 2021, doi: <https://doi.org/10.31289/jite.v5i1.5308>.
- [18] S. E. Sitepu, U. Satya, and T. Bhinneka, “Low-Resource Single-Domain Machine Translation untuk Bahasa Karo-Indonesia Pendahuluan,” vol. 1, no. 4, pp. 59–66, 2023. doi: <https://doi.org/10.31004/ijme.v1i4.21>
- [19] D. Torregrosa *et al.*, “Leveraging Rule-Based Machine Translation Knowledge for Under-Resourced Neural Machine Translation Models,” *Proc. Mach. Transl. Summit XVII Transl. Proj. User Tracks*, vol. 2, pp. 125–133, 2019, [Online]. Available: <https://aclanthology.org/W19-6725.pdf>.
- [20] J. Daems and L. Macken, “Interactive adaptive SMT versus interactive adaptive NMT: a user experience evaluation,” *Mach. Transl.*, vol. 33, no. 1–2, pp. 117–134, 2019, doi: 10.1007/s10590-019-09230-z.
- [21] Y. Liu, J. Gu, N. Goyal, X. Li, and S. Edunov, “Multilingual Denoising Pre-training for Neural Machine Translation,” vol. 8, pp. 726–742, 2020.
- [22] Z. Yu, Z. Yu, J. Guo, Y. Huang, and Y. Wen, “Efficient Low-Resource Neural Machine Translation with,” vol. 19, no. 3, pp. 1–13, 2020.
- [23] D. Puspitaningrum, “A Study of English-Indonesian Neural Machine Translation with Attention (Seq2Seq, ConvSeq2Seq, RNN, and MHA): A Comparative Study of NMT on English-Indonesian,” *ACM Int. Conf. Proceeding Ser.*, pp. 271–280, 2021, doi: 10.1145/3479645.3479703.
- [24] D. Datta, P. E. David, D. Mittal, and A. Jain, “Neural Machine Translation using Recurrent Neural Network,” *Int. J. Eng. Adv. Technol.*, vol. 9, no. 4, pp. 1395–1400, 2020, doi: 10.35940/ijeat.d7637.049420.
- [25] S. Iida, R. Kimura, H. Cui, P.-H. Hung, T. Utsuro, and M. Nagata, “A Multi-Hop Attention for RNN based Neural Machine Translation,” *Proc. 8th Work. Pat. Sci. Lit. Transl.*, vol. 2018, pp. 24–31, 2019, [Online]. Available: <https://aclanthology.org/W19-7203>.
- [26] A. Shewalkar, D. nyavanandi, and S. A. Ludwig, “Performance Evaluation of Deep neural networks Applied to Speech Recognition: Rnn, LSTM and GRU,” *J. Artif. Intell. Soft Comput. Res.*, vol. 9, no. 4, pp. 235–245, 2019, doi: 10.2478/jaiscr-2019-0006.
- [27] Y. Dong, “RNN Neural Network Model for Chinese-Korean Translation Learning,” *Secur. Commun. Networks*, vol. 2022, 2022, doi: 10.1155/2022/6848847.
- [28] R. Achmad, Y. Tokoro, J. Haurissa, and A. Wijanarko, “Recurrent Neural Network-Gated Recurrent Unit for Indonesia-Sentani Papua Machine Translation,” *J. Inf. Syst. Informatics*, vol. 5, no. 4, pp. 1449–1460, 2023, doi: 10.51519/journalisi.v5i4.597.
- [29] M. Wahyuni, H. Sujaini, and H. Muhandi, “Pengaruh Kuantitas Korpus Monolingual Terhadap Akurasi Mesin Penerjemah Statistik,” *J. Sist. dan Teknol. Inf.*, vol. 7, no. 1, pp. 20–26, 2019, doi: <https://dx.doi.org/10.26418/justin.v7i1.27241>.
- [30] B. Premjith, M. A. Kumar, and K. P. Soman, “Neural machine translation system for English to Indian language translation using MTIL parallel corpus,” *J. Intell. Syst.*, vol. 28, no. 3, pp. 387–398, 2019, doi: 10.1515/jisys-2019-2510.

Title of manuscript is short and clear, implies research results (First Author)