Algorithm Design for Converting Indonesian Latin to Sasak Latin Using the Sequence-To-Sequence Transformers Method

Muhammad Giri Restu Adjie^{[1]*}, Ramaditia Dwiyansaputra^[1], Fitri Bimantoro^[1], Arik Aranta^[2] ^[1]Dept Informatics Engineering, Mataram University Mataram, Lombok NTB, INDONESIA

> ^[2]Departmen of Informatics, Sepuluh Nopember Institute of Technology Surabaya, East Java, INDONESIA

Email: giriadjie44@gmail.com, rama@unram.ac.id, bimo@unram.ac.id, 7025231011@student.its.ac.id

*Corresponding Author

Sasak is a regional language spoken in West Nusa Tenggara, facing challenges in maintaining its use, especially among the younger generation, due to the dominance of Indonesian in formal settings. This research explores these challenges and highlights the importance of preserving the Sasak language as an important social and cultural identity for the Sasak people. This research aims to develop a machine translation system for converting Indonesian to Sasak using the Sequence-to-Sequence Transformer method. By using the Transformer model with its encoder-decoder based architecture, this research translates Indonesian text into Sasak, by utilizing the Rule-Based method for preprocessing the dataset. The dataset used consists of more than 85,290 lines of Indonesian-Sasak text pairs, which are divided into training, validation and testing sets. The training carried out by this model achieved a final result of accuracy after 30 epochs of 0.99 and validation accuracy of 0.98 and with a score of 0.6075 in the Bilingual Evaluation Understudy (BLEU) evaluation. This shows the model's strong ability to produce accurate translations, even though the Sasak language is complex. This research is not only to preserve the Sasak language but also opens new avenues for future researchers in language processing and preservation, especially for languages with less resources such as Sasak.

Key words: Sequence-to-sequence Transformer, Sasak Language, Translation Machine, Keras, Tensorflow.

I. INTRODUCTION

Language is a means of communication composed of units such as words, phrases, clauses, and sentences, conveyed either orally or in written form. This is only one of many definitions of language. For instance, language can also be described as a human communication system that organizes sounds or written symbols into structured forms, creating larger elements such as morphemes, words, and sentences. Around the world, thousands of languages exist, each with its distinct system known as grammar, such as the grammar of Indonesian, English, Japanese, and many others[1]. Indonesia has around 733 local languages that need to be preserved, many of which are in critical status due to the decreasing number of native speakers who no longer use or pass on the language to the next generation. Over time, the number of these languages continues to decline and is threatened with extinction, some of which are even heading towards extinction[2].

The Sasak language is one of the regional languages spoken in West Nusa Tenggara. This language is used as a medium of communication by the Sasak people who live on Lombok Island. The Sasak people are the indigenous inhabitants of Lombok Island, forming the majority of the population. The Sasak language is still used as a medium of communication among the Sasak people. However, from observations, it is very rare for the Sasak people to use the Sasak language even in formal gatherings, although only Sasak ethnic groups are present. Even in Sasak community meetings, Sasak community leaders often deliver speeches in Indonesian. In government or community organization meetings, even when attended solely by Sasak ethnic groups, communication is frequently conducted in Indonesian[3]. The Sasak language represents the identity of the Sasak people and plays a vital role in their social and cultural life. However, with the rapid development of modern times, there is concern about the decreasing usage of the Sasak language, especially among the younger generation, making it increasingly rare and rigid, thus threatening the continuity of their culture in the digital era.

Sasak is also a language with limited resources unlike languages with abundant resources such as English, French and German. According to research conducted in 2023 on Kurdish language translation with the transformer method which is considered one of the languages with limited resources and produced a score of 0.45 on the Bilingual Evaluation Understudy (BLEU) evaluation which indicates high translation quality according to BLEU standards[4]. Based on this concern, researcher is interested in conducting a study titled "Design and Development of an Algorithm for Converting Indonesian into Sasak Language Using the Sequence-To-Sequence Transformer Method." This research utilizes the Transformer method as a translator from Indonesian to Sasak. Transformers are the first transduction model that relies entirely on self-attention to compute input and output representations. The Transformer is a structure based on an encoder-decoder model, where the encoder consists of multiple layers that process input iteratively one by one, while the decoder consists of a set of encoding layers that do the same with the encoder output[5]. This research is expected to support the development of a better and more relevant translation engine for Sasak language, which is categorized as a language with limited resources. In addition, this research also aims to help people who want to learn Sasak as well as contribute to preservation efforts through documenting the language.

II. LITERATURE REVIEW

Research on the application of the Transformer method and language translation has been conducted by several previous researchers. The following are some studies that can serve as references for this research.

The first study is about conversion Indonesian voice into sasak latin using Dictionary Based method. In this research, a system was designed that implements Google's speech-to-text API to translate Indonesian words or sentences into the Sasak language. The testing process involved translating 25 sentences taken from the Sasak-Indonesian dictionary, consisting of 117 words. There were two stages of testing in this study. The first test aimed to determine the accuracy of the translation from Indonesian to Sasak using the dictionary method. The second test aimed to assess the accuracy of the Google Speech API in recognizing voice input and converting it into text. From the first test, the system's accuracy in translating Indonesian into Sasak using the dictionarybased method was 100%, with a 0% error rate. Meanwhile, in the second test, the system successfully implemented the Google Speech API to translate Indonesian words or sentences into Sasak with an accuracy rate of 99.14%[6].

The second study is titled "A Transformer-based Neural Network Machine Translation Model for the Kurdish Sorani Dialect." The Transformer model is one of the latest innovations in text translation development, utilizing attention mechanisms to outperform previous models such as sequence-to-sequence in terms of performance. This model has proven to be highly effective in resource-rich languages like English, French, and German. In this study, the researchers explored the application of a modified Transformer architecture to train a model on low-resource languages, such as Kurdish. The study introduced the first Transformer-based neural machine translation model for the Kurdish language, utilizing shared vocabulary units across the parallel corpus. To achieve this, all available Kurdish-English parallel corpora were combined to form a larger dataset, which was then used to train the proposed Transformer model. The experimental results showed that this model was capable of delivering good performance in translating Kurdish texts, with a score of 0.45 in the Bilingual Evaluation Understudy (BLEU) evaluation, indicating high translation quality according to BLEU standards[4].

The third study is titled "A French-to-English Machine Translation Model Using Transformer Network." Traditional machine translation based on Recurrent Neural Networks (RNN) has two major drawbacks: first, the translation process is done word by word, resulting in slow training speed; second, in translating long sentences, RNNs are prone to gradient vanishing and gradient explosion problems, which ultimately lower the translation accuracy. To address these weaknesses, this paper introduces a machine translation design based on the Transformer architecture, implemented using PyTorch. Unlike RNN models, Transformers use an attention mechanism to overcome the limitations of RNNs. addressing efficiency effectively issues and ineffectiveness with long sequences. In this study, the Transformer-based French-to-English translation system successfully achieved an 80% accuracy rate in translating from French to English after practical training and application[7].

The fourth study is titled "Robust Neural Language Translation Model Formulation using Seq2seq approach" This research utilizes a method that sequences a robust model with proven success in language translation tasks involving encoding and re-encoding. The approach uses a language transformer model based on a sequence-tosequence framework, utilizing a Long Short-Term Memory (LSTM) network to transform the input sequence into a fixed-dimensional vector. Subsequently, another deep LSTM generates the target sequence from this vector. The efficiency of the model was assessed using BLEU scores, which showed that the BLEU scores of LSTMs were negatively affected by words that were not in the vocabulary. However, LSTM handles sentences of varying length effectively. This study shows that the deep configuration of LSTM for English to Japanese translation achieves much faster performance on both GPU and CPU. Diverse datasets were combined to assess the robustness of the model through BLEU scores. In the end, combining two different datasets resulted in a higher BLEU score of 0.401, which marked the best performance achieved in this study[8].

The fifth study is titled "Analisis Pengembangan Model Neural Machine Translation (NMT) dengan Transformer untuk Penerjemah Bahasa Indonesia ke Bahasa Gayo " This study adopted a training approach using a parallel corpus collected from the Indonesian-Gayo II Dictionary. The dataset underwent a series of preprocessing steps to prepare it for use in model training. The pre-processing steps included case folding and removing punctuation. The dataset was divided into three parts consisting of training data, validation data, and test data. The researchers also performed data augmentation using the MixSeq method to increase the diversity and size of the training data. Using an optimized Transformer architecture, experimental results showed an accuracy rate of 90% for the augmented training data. The evaluation using the BLEU score resulted in a score of 79.90[9].

The sixth study entitled "Mesin Penterjemah Bahasa Indonesia-Bahasa Sunda Menggunakan Recurrent Neural Networks" This study builds an Indonesian to Sundanese translator. The stages used start from pre-processing using text preprocessing and word embedding Word2Vec and the approach used is Neural Machine Translation (NMT) with Encoder-Decoder architecture in which there is a Recurrent Neural Network (RNN). Testing on the research resulted in an optimal value by GRU of 99.17%. The model using Attention gets 99.94%. The use of the optimization model gets optimal results by Adam 99.35% and the results of BLEU Score with optimal bleu 92.63% and brievity penalty 0.929. The results of the translation machine produce training predictions from Indonesian to Sundanese when the input sentence matches the corpus and the translation results are less appropriate when the input sentence is different from the corpus[10].

A. Sasak Tribe

The Sasak people, as speakers of the Sasak language, are socially divided into two groups: the noble class (menak) and the common people (jajarkarang/bulu ketujur). The noble class is further divided into two levels, while the common people consist of only one level, as affirmed by Mahyuni: "Traditionally, Sasak people were divided into four social classes: Raden 'prominent nobles', menak and perwangse 'ordinary nobles', and jajarkarang or bulu ketujur 'commoners'."[3].

B. Sasak Language

The Sasak language is one of the regional languages in West Nusa Tenggara. This language is used as a means of communication by the Sasak people who reside on the island of Lombok. The Sasak people are the native inhabitants of Lombok Island and make up the majority of the population. The Sasak language is still used as a medium of communication among the Sasak people[3].

C. Keras

With the growing development and research in Deep Learning, many libraries have emerged focusing on artificial neural networks. One such example is Keras. Keras is a high-level neural network library written in Python and capable of running on top of TensorFlow, CNTK, or Theano[11].

D. Tensorflow

The development of the Deep Learning field has been facilitated by the abundance of libraries and Application Programming Interfaces (APIs). The library used is TensorFlow, which serves as an interface for expressing machine learning algorithms and executing commands using the information it has about the objects or recognized targets, as well as differentiating one object from another. TensorFlow features the ability to run model training using a Central Processing Unit (CPU) and a Graphics Processing Unit (GPU). However, in this implementation, model training will be conducted using the CPU feature[12].

E. Translation Machines

Machine translation is a technique for converting source sentences from one natural language to another using a computerized system, eliminating the need for human intervention[13].

F. Rule-based

Rule-Based Systems are software frameworks that apply expert knowledge represented as rules in a specific domain to address problems. These systems are straightforward and can be adapted for various problems. However, as the number of rules increases, maintaining the system becomes more challenging, and errors can arise. Rule-based methods operate by using a rule base that summarizes all relevant knowledge about the problem, encoded as if-then statements containing data, statements, and initial conditions. The system evaluates the rules and executes the corresponding later conditions when a match is found. This process continues iteratively until one of two outcomes is achieved: a matching rule is identified, or the system exits the loop (terminates) if no suitable rule is found[14]. Apart from rule-based relying heavily on the quality and scope of human-defined rules, rule-based is not suitable for training with large corpus datasets. Rule-based is also weak in handling long sentence contexts or complex idioms[4].

G. Deep Learning

Deep Learning is a learning approach that uses a multilayered artificial neural network modeled after the human brain, with interconnected neurons forming an intricate network. Also referred to as deep structured learning, hierarchical learning, or deep neural learning, it uses multiple nonlinear transformations to process data. Deep Learning can be understood as the convergence of machine learning and artificial neural network techniques, which enables advanced computational capabilities[15].

H. Transformer

The Transformer method is a deep learning model that was first published in the paper "Attention is All You Need" in 2017. Transformers are the first transduction models that rely entirely on self-attention to compute the representation of their input and output. Transformers have an encoderdecoder structure, where the encoder consists of a combination of encoding layers that process the input iteratively, one by one, while the decoder consists of a set of encoding layers that perform the same operation on the encoder's output[5]. A more in-depth abstract representation is obtained through the use of background representations (embeddings) generated directly from the data, allowing the model to capture complex relationships between words and phrases through self-attention mechanisms. This approach enables learning directly from the data without the need for explicit rule programming, thus increasing the flexibility and generalization capabilities of the model. In addition, the Transformer architecture has demonstrated its ability to achieve higher BLEU scores on English-German and English-French

translation tasks compared to RNN and LSTM-based models[5].

III. METODOLOGY

A. Tools and Materials

The tools used in this research consist of hardware and software. The materials used include datasets and several literature sources found by the researchers. The tools and materials utilized in this study are detailed in Table I and Table II.

A.1.Hardware

TABLE I. TABLE OF HARDWARE

No	Hardware	Specification
1	Laptop	Acer Aspire E5-476G dengan Processor Intel(R) Core(TM) i5-8250U CPU @1.60GHz 1.80GHz, RAM 12,0 GB

A.2. Software

TABLE II. TABLE OF SOFTWARE

No	Hardware	Specification
1	Operating System	Windows 11 64-bit
2	Programming Language	Python
3	Code Editor	Jupyter Notebook



Fig. 1. Figure of research flow

A.3. Research Materials

The materials used include a dataset consisting of English vocabulary obtained from Anki, which is then translated into Indonesian using Google Translate. Subsequently, the Indonesian vocabulary acquired is translated into Sasak, with references taken from the Sasak dictionary published by the Nusa Tenggara Barat Language Office.

B. Research Flow

In general, the flow of the research conducted in this study can be illustrated as a flowchart shown in Figure 1. Before the design and development of the system, a literature review was conducted to enhance the knowledge and understanding of what will be researched next. The literature review can be carried out by reading various references related to similar research, such as journals, articles, and books.

C. System Design

Figure 2 is an overview of the system design of the research methodology that will be applied for the conversion of Indonesian Latin into Sasak Latin.



Fig. 2. Figure of system design

E-ISSN:<mark>2541-0806</mark> P-ISSN:<mark>2540-8895</mark>

A.1. Preparing Dataset

Download Dataset

Anki is an application used to memorize information through a system of reminder cards or "flashcards". Since Anki is publicly available, this research uses the dataset downloaded from Anki, which is English-Spanish text data from Tatoeba Project (translation from English to Spanish data available at <u>http://www.manythings.org/anki</u>).

• Filter Dataset

After the data in the form of text is downloaded where English and Spanish in the text are distinguished using indents (tab). Furthermore, the Spanish language is removed using a program made using the python language with the aim of leaving English so that it can then be translated into Indonesian because Spanish has a more complex word structure that will significantly affect the results of translation into Indonesian.

- Translate English into Indonesia To translate English to Indonesian, the googletrans library in Python is used. Googletrans is a Python library that utilizes google translate and does not require API configuration so it is suitable for use in the English-Indonesian translation process[6].
- Filter Dataset

The dataset that has been translated into English-Indonesian is then removed from the English language using a program created using Python like the removal of Spanish.

Remove Redundant Dataset

Indonesian language datasets still have a lot of repetition of the same text or sentence. Therefore, to obtain a good dataset, redundant text or sentences are removed from the dataset using a Python program.

Sorting Dataset

Sort the Indonesian language dataset by character length so that the dataset is neatly arranged from the shortest to the longest character.

Collect Sasak Dataset

To obtain the Sasak language corpus, the dataset that has been translated into Indonesian will then be translated into Sasak using the Rule-based method due to its simple operation and then save the results in a file in txt format. According to research conducted in 2024 on the conversion of Indonesian voice into Sasak language Latin text using the Rule-based method achieved an accuracy rate of 99.14% in text-to-text translation[6]. Because the Rule-based method uses "if-then" rules, it requires an additional dataset whose content is Indonesian text with its Sasak language translation word-for-word. Therefore, the translation of Indonesian-Sasak vocabulary was carried out using the "Kamus Sasak Indonesia" published by the West Nusa Tenggara Office, Language Development and Development Agency, and the Ministry of Education and Culture in 2017 as a reference. However, the writing of this dictionary still uses a phonetic system for the division of dialects in its writing, meaning that

the dataset taken from the dictionary still uses mixed dialects.

Combine Dataset

Combine the Indonesian language dataset with the Sasak language so that it becomes one file that will be used as an input-target pairs dataset with Indonesian as input and Sasak language as the target with the file name "SasakId.txt", generated 85,290 lines and 6,282,484 characters as shown in table III.

TABLE III.	TABLE OF DATASETS AMOUNT
	THEE OF ENTITED TO THE OTHER

No	Dataset	Total Dataset
1	Line	85.290
2	Character	6.282.484

A.2. Split Dataset

Before the dataset undergoes preprocessing, it will be divided into three parts: the training set, test set, and validation set. The training set is the data used for training, where the model will learn from the samples contained within it. The validation set is the data used to evaluate the model during the training process and provides an indication of the model's performance on data that it has not seen before. Meanwhile, the test set is the data used for testing. Each sample in the test set is a sample that the model has never encountered during training, making the test set useful for objectively measuring the model's performance.

A.3. Preprocessing

Preprocessing is the process of altering the structure of text to meet specific needs. The data processed consists only of the training sets, and the preprocessing methods applied are:

TextVectorization for processing text in two languages, namely Indonesian and Sasak. Two layers of TextVectorization are used, one for Indonesian and the other for Sasak. The function of the TextVectorization layer is to convert raw text into numerical forms represented by sequences of numbers or integer sequences. Each number in the sequence represents the position or index of a word in the vocabulary or list of words recognized by the model. The Indonesian layer will use default string standardization and split scheme, while the Sasak layer will utilize custom standardization and split scheme. The default string standardization used in the Indonesian layer involves removing punctuation marks (such as periods, commas, exclamation marks) from the text, leaving only the words. The custom standardization applied in the Sasak layer is similar to that used in the Indonesian layer, with the modification of the standardization process by employing lower casing, which converts all words to lowercase. The split scheme involves separating the text based on spaces after punctuation marks have been removed, meaning that every word separated by

spaces will be considered a separate unit in the vectorization process.

• Dataset format: the dataset is structured in the form of pairs (inputs, targets) that will be used for training. Inputs refer to the model's input in the form of a dictionary with two keys: encoder_inputs and decoder_inputs. Encoder_inputs are the source sentences that have been vectorized into sequences of numbers, representing the original text to be translated or processed by the model. Decoder_inputs are the target sentences processed thus far, which are the words in the target sentence from position 0 to N, used by the model to predict the next word (words N+1, N+2, etc.) in the target sentence. Targets are the target sentence contains the words that the model will attempt to predict in the next step.

A.4. Training

The Transformer is a deep learning model with an encoder-decoder architecture[5]. The transformer model training conducted in this study used[16]:

Encoder. The encoder consists of NNN layers, where each layer contains a multi-head attention sublayer and a feed-forward network. Each sublayer applies a residual connection and layer normalization, *LayerNorm* (x + Sublayer(x)), where x is the embedding vector and *Sublayer(x)* is the function implemented by the sublayer.

Decoder. The decoder consists of N layers. The decoder has three sublayers, two of which are the same as the sublayers in the encoder, while the other sublayer is a masked multi-head attention sublayer.

Embedding. The transformer converts the token into a vector with dimension d_{model} . The weight matrix in the embedding layer is multiplied by $\sqrt{d_{model}}$.

Positional Embedding. Due to its non-recurrent nature, Transformer uses positional embedding to provide information about the position of tokens so that the model can understand the word order. Each word or token at a position is represented with a unique embedding vector determined by its relative position (i-th position). For each token at position i, positional embedding provides a vector representing that position in the sequence of words which can help the model consider the relationship between words based on their position.

Self-Attention. Self-attention is a mechanism to calculate the connectedness between words in a sentence. To calculate the attention value, the input vector x is projected into the query (Q), key (K), and value (V). Query and key have dimension d_k , while value has dimension d_v .

$$Attention(Q, K, V) = softmax\left(\frac{QK^{T}}{\sqrt{d_{k}}}\right)$$
(1)

Multi-Head Attention. Instead of calculating one attention value with d_k equal to d_{model} , multi-head attention allows the attention value to be computed *h* times with $d_k = d_n = d_{model}/h$.

$$MultiHead(Q, K, V) = Concat(head_i, ..., head_h)W^o$$
(2)

$$head_{i} = Attention(QW_{i}^{Q}, KW_{i}^{K}, VW_{i}^{V})$$
(3)

Where is the weight matrix $W_i^Q \in R^{d_{mocel} \times d_k}$, $W_i^K \in R^{d_{mocel} \times d_k}$, $W_i^V \in R^{d_{mocel} \times d_k}$, dan $W_i^O \in R^{hd_v \times d_{model}}$.

Feed-Forward Network. Each layer in the encoder and decoder has a feed-forward network consisting of two linear transformations with a ReLU activation function between them. The dimensions of the inputs and outputs are d_{model} , while the inner-layer has the dimensions d_{ff} .

$$FFN(x) = max(0, xW_1 + b_1)W_2 + W_2$$
(4)

A.5. Testing

The testing process is carried out by providing the model with unseen test data, which has never been used during training, to evaluate how well the translation model performs. The test pairs consist of sentences in Indonesian, which the model will then translate into Sasak. This step aims to assess whether the model is capable of generating accurate translations from the trained data.

$$Accruation = \left(\frac{N_{correct}}{N_{Total}}\right) \times 100\%$$
(5)

A.6. BLEU Evaluation

BLEU score is a string matching-based algorithm that provides a basic metric to measure the quality of the output which uses candidate and reference to assess the score.

$$BLEU = BP \cdot \exp(\sum_{n=1}^{N} \omega_n \cdot \log \rho_n)$$
(6)

IV. RESULT AND DISCUSSION

A. Training Configuration

In this study, the dataset consists of fragments of words or complete sentences in two different languages: Indonesian and Sasak. A total of 85.290 lines of word fragments or sentences are used, stored in a TXT format. During the training phase, the dataset includes data in both Indonesian and Sasak, which are separated by indents (tab). Before use, the dataset is divided into three parts: train pairs, validation pairs, and test pairs. Train pairs are used for training, allowing the model to learn from the samples within them; validation pairs are used to evaluate the model during the training process; and test pairs are used for testing. The training is conducted using a composition of 59,703 training pairs, 12,793 validation pairs, and 12,793 test pairs with ratio 70/15/15 as shown in table IV.

TABLE IV. TABLE OF SPLITTED DATASETS

No	Datasets	Jumlah
1	Train pair	59703
2	Validation pair	12793
3	Test pair	12793
Total		85289

The model is trained for 30 epochs, batch size 64, and acrivation function used softmax with optimizer used is RMSProp (Root Mean Square Propagation) with the default configuration of a learning rate 0.001, $\rho = 0.9$, and $\epsilon = 1e-07$. the details can be seen in the Table V.

No	Datasets	Jumlah
1	Epoch	30
2	Batch size	64
3	Activation function	Softmax
4	Optimizer	RMSProp (default, 0,001, $\rho = 0.9$, and $\epsilon = 1e-07$)

TABLE V. TABLE OF HYPERPARAMETER CONFIGURATION

The hyperparameter configuration shown in table V was used to provide sufficient training time, maintain a balance between the efficiency and generalization ability of the model, and ensure stability and optimal performance in the training process. Training for 30 epochs was chosen to allow sufficient time for the model to learn from the data and achieve convergence, while minimizing the risk of overfitting. A batch size of 64 was chosen as this size often provides a good balance between training efficiency and the model's ability to generalize. A softmax activation function was applied to the output layer as it is suitable for multi-class classification tasks, allowing the model to generate probabilities for each class. RMSProp was chosen as the optimizer due to its ability to dynamically adjust the learning rate during training, which is helpful in handling fluctuating gradients, especially in complex transformer models. The learning rate of 0.001, which is a commonly used default value, provides stability and consistent results in many experiments, ensuring the training process runs at the right speed.

B. Training Data

Layer (type)	Output Shape	Param #	Connected to
Encoder_input (In putLayer)	(None, None)	0	-
Positional_embed ding_2 (PositionalEmbed ding)	(None, None, 256)	3,845 ,120	Encoder_input s[0][0]
decoder_inputs (InputLayer)	(None, None)	0	-
<pre>transformer_enco der_1 (TransformerEnco der)</pre>	(None, None, 256)	3,155 ,456	Positional_em bedding_2[0][0]
<pre>functional_11 (Functional)</pre>	(None, None, 15000)	12,95 9,640	Decoder_input s[0][0], transformer_e ncoder_1[0][0]
Total params: 19 960 216 (76 14 MB)			

 $TABLE \ VI. \ \ TABLE \ OF \ TRANSFORMER \ SUMMARY \ MODEL$

Trainable params: 19,960,216 (76.14 MB) Non-trainable params: 19,960,216 (76.14 MB)

Model summary shown as Table VI used to describe in detail the main components of the Transformer model used, namely the type of layer, output shapes, the number of parameters involved, and the layers that are interconnected in the model. The data flows in the model:

- Encoder input Sentences from the source language or Bahasa Indonesia enter the Encoder_input layer and then the source token is converted into embedding with the addition of position information by Positional_embedding_2.
- Processing by Encoder The token representation that has been processed by Positional_embedding_2 goes into the transformer_encoder_1, where each token is enriched with context using self-attention and feed-forward layers.
- Decoder input

Sentences from the target language or Sasak language are fed into the decoder_inputs and during the decoding process the target tokens are processed with the representation of the encoder in the functional_11 layer, resulting in a probability distribution for each token in the target vocabulary.

Output

The functional_11 layer outputs the probability distribution for each token (1500 tokens) in the target vocabulary or Sasak vocabulary. This process is used to predict the next token in the translation sequence.







Fig. 4. Figure of training and validation loss graphic

During the training process using the training pairs and validation pairs, the results obtained in the first epoch indicate a relatively high accuracy of 0.86 and a loss of 1.37, with a considerable gap between the validation accuracy of 0.96 and validation loss 0.31. After 10 epoch the model began to adjust its weight better, and accuracy began to improve with accuracy of 0.98 and loss of 0.07. and validation accuracy of 0.98 and validation loss of 0.14. Convergence is achieved when the accuracy of the model stabilizes and no longer increases significantly after a few epochs. This indicates that the model has learned a good representation and is ready to make predictions with a high level of accuracy. However, after 30 epochs, the results improved significantly, achieving an accuracy of 0.99, a loss of 0.02, and a validation accuracy of 0.98 and validation loss of 0.15, as shown in figure 3 and figure 4.

C. Result

In the results section, evaluation will be carried out using two scenarios. The first uses test pairs, namely datasets that have never been seen by the model to see the translation results after the model has been trained for 30 epochs. The second model will be evaluated using the Bilingual Evaluatin Study (BLEU) to assess the BLEU score obtained after the model has been trained for 30 epochs.

A.1. Decoding test sentences

In the decoding test sentences section, random selections of previously unseen English sentences from the test data are translated using the trained transformer model. The process begins with converting the input sentence into tokens and then proceeding with the translation step by step. At each step, the model predicts the next word based on the sentences translated so far and the provided input sentence. This continues, with the model choosing the word that has the highest probability at each step, until the sentence is complete or reaches a set length limit. This approach enables the model to generate precise translations that align with the context of the input sentence.

TABLE VII.	TABLE OF PREDICTION
------------	---------------------

Test Pairs	Prediction
ini akan sangat mahal.	niki gen sangat mahel
pria itu memiliki sesuatu di bawah mantelnya.	mame nike bedoe sesuatu leq bawaq mantelne
tom menahan senyum.	tom menahan kemos
label harga masih ada di kemeja yang dikenakan tom.	label aji masi araq leq ojokmeje siq dikenakan tom
dia jarang pergi ke bioskop.	dia jarang pergi ke bioskop.
kami akan memeriksanya sekarang.	kami gen memeriksane nengke
tom membuat daftar lagu yang tidak dia sukai.	tom mempiaq daftar agol siq endeq niye demeni
saya tidak mengingat anda.	tiang endeq mengingat pelungguh
maaf saya mengganggu anda.	maaf tiang mengganggu pelungguh
aku bisa menciummu sekarang.	aku bau menciumde nengke

Table VII shows the model translation results after 30 epochs. While decoding test sentences gives an idea of the model's ability to translate sentences, the resulting translations are not sufficient to thoroughly evaluate the quality of the model. This is because the decoding process only measures the model's ability to produce sentences sequentially, without considering the semantic match and overall sentence structure. Therefore, to obtain a more comprehensive and objective evaluation, an additional evaluation using the BLEU (Bilingual Evaluation Understudy) metric is required, which can measure the extent to which the translation produced by the model matches the human translation based on the n-grams embedded in the target sentence.

A.2. Bleu Evaluation

BLEU (Bilingual Evaluation Understudy) evaluation is used to measure the quality of the translation produced by the model by comparing it against the provided reference translation. In this evaluation, candidates refer to the translation produced by the model, while references are human translations or validated translations that serve as a reference, the reference used is the Sasak language that has been published in the form of Sasak folklore. BLEU calculates a score based on the n-gram similarity between candidates and references, taking into account higher ngram precision as well as penalties for sentences that are too short. A higher BLEU score indicates that the model translation is more similar to the reference translation, indicating better translation quality.

$TABLE \ VIII. \ \ TABLE \ OF \ CANDIDATES \ AND \ References$

Candidates tetu-tetu santer bodo, sesuai kance aranne loq sesekeq. sekalipun ngeno, sesekeq santer tesayang kangen siq inakne. sopoq jelo, inakne suruq sesekeq lalo beli kemek ojok peken, lalo begawean jari saudagar kemek, saudagar saq jujur kance rajin.

References

teceritaan, leq zaman laeq araq sopoq cerite. sepasang senine sememe sanget patuh. leq akhir hayatne simeme bepesen tipak seninakne si nyeke betian. "lamun lahir anakte meme, ndak lupak beng ye aran loq sesekeq. kontek cerite lahir meme, mukne teparan log sesekeg. sesekeq artine bodo. sopok jelo leq waktu genne araq subuh, manuk ngungkung tende benar. selapuk kemanukan tarik muni, tesambut isik sueren bang leg masjid. log sesekek tures lalo sembahvang. selese sembahyang, sesekeq lalo tulung inakne nyepu meriri, ronas piring, mopoq natap dait siram tetaletan. tetu- tetu mule bodo, sesuai isik aranne Lok sesekeq. sekalipun sak ngeno, sesekeq sanget isine tetunah kangen isik inakne. sopoq jelo, inakne suruk sesekeq lalo beli kemek ojok peken, lalo nguruk jeri saudagar kemek, saudagar si jujur dait pecu. sesekeq ndek neuah tolak perintah inakne. sambil terenyuk lalok si tesuruk beli kemek. si pikiranane berembe entan yak jauk kemek. dateng leg peken, langsung ye pileq kemek si ndek boke atau ntek. yahne beli kemek sino, bingiung ye loq sesekeq. berembe bae ntan yak jauk kemek sino. " oh, eku bingung jauk kemek sine. bedagang doang dekke tao, epe legi jauk *kemek. oh, eku dait akal... mem ! " unin sesekeq. banjur* boyakne telu, beterusne totos kemek sino ntan-ntan sopoq. Suahan sino boyakne teli, beterus perentokne kemek sino jeri sekeq, langsung teoros ojok sopok taok. Sesekeq lampaq ndek nearak kereguan. Sesekeq pencar kemekne teapek yakne bedagang jeri saudagar kemek. Ndek arak bae dengan mele beregak sekek-sekek. Bahkan lueq dengan si bengak lalok gitak kemek tarik tepong. Ndek jak arak leku kemekne, Sesekeq ulek lalo ngelapur tipak inakne. Muni inakne "O, gamak anakku, tetu jak luek dengan dateng, tetujak luek pemborong, laguk pasti pede bengak si gitak kemekde si selapukne tepong". Lamun ngeno, becat de lalo malik ojok peken. Beliang te bebek, laguk ndak bae beli kelinci laun salak belinde....etc.

Using the candidates and references as shown in Table VIII for the BLEU evaluation resulted in a score of 0.6075 which is considered a fairly high result. The comparison of the number of candidates and references is classified using the number or length of words. A large number of references were used with the aim of improving the accuracy and representativeness assessment of translation quality. The comparison for candidates and references is about 1800 words for references and 37 words for candidates. The word count comparison diagram of candidates and references can be seen in figure 5.



Fig. 5. Figure of Comparison of Candidates and References

Table IX shown as classification of translation quality based on BLEU Score to determine the evaluation result of the Transformer training model.

TABLE IX. TABLE OF CLASSIFICATION BLEU SCORE

BLEU Score Range	Quality
0 - 0.10	Very Poor
0.10 - 0.20	Poor
0.20 - 0.30	Fair
0.30 - 0.40	Good
0.40 - 0.50	Very Good
0.50 - 1.00	Excellent

Based on the classification shown in table IX, The evaluation results using BLEU, which reached a score of 0.6075, show that the translation produced by the model has a moderate level of similarity with the reference translation. While this score indicates that the model is able to produce translations that are quite relevant to the target sentence, it also shows that there is still room for improvement, both in terms of n-gram precision and in the model's ability to capture more complex contexts. Therefore, while the BLEU results reflect the model's good performance, further research is needed to improve the accuracy and quality of the translations.

V. CONCLUSION AND SUGGESTION

Based on the test results, the Transformer model shows excellent performance in translating Sasak language, with accuracy reaching 0.99 after 30 training epochs, loss of 0.02, validation accuracy of 0.98, and validation loss of 0.15. The BLEU score of 0.6075 also showed positive results. This shows that the model can produce translations that are quite relevant to the target sentence, but there is still room for improvement. The model successfully proved that Neural Machine Translation (NMT) technology can be well applied to Sasak language, despite resource limitations, especially in terms of data amount and dialect diversity. Sasak has many dialects that can affect translation quality, so the model needs to be further tested with a corpus that includes a variety of dialects to test its ability to handle more complex language variations.

improve performance, future research is То recommended to develop a larger and more diverse corpus, and use more comprehensive evaluation techniques to measure translation quality. Given the diversity of Sasak language dialects, the model needs to be adapted to specific approaches for these dialects, one of which is by using a lock tokenization approach. This approach can help reduce the problem of mismatch between dialectal variations and translation models by suggesting a more efficient and relevant division of tokens within each dialect. The findings of this study open up opportunities for researchers to develop the Transformer model with parameter adjustments for Sasak language, as well as creating a larger and more structured parallel corpus for the model to be applied more effectively in language translation applications.

REFERENCES

- [1] R. Brianto, M. Pakaja, C. Ishak, Musrifa, M. Hunowu, and A. A. Mohamad, "Novateur Publication, India Proceedings of International Seminar on Indonesian Lecturer is Born to Report Regularly Communicative Language Therapy for Alzheimer's Dementia in Neurolinguistics," in *International Seminar on Indonesian Lecturer is Born to Report Redularly*, India: Novateur Publication, May 2023, pp. 466–472.
- [2] A. Rahima and U. Batanghari Jambi, "Pengabdian Deli Sumatera Revitalisasi Bahasa Daerah Hampir Punah Sebagai Dokumentasi Bahasa," *Journal Law of Deli Sumatera*, vol. 3, no. 1, pp. 56–61, Jul. 2024, [Online]. Available: http://118.98.223.79/petabahasa/.
- [3] Mugni, "Pemertahanan Bahasa Sasak pada Keluarga Bangsawan Lombok (Studi Etnografi di Kabupaten Lombok Timur)," *Jurnal Linguistik, Sastra, dan Pendidikan (JURNALISTRENDI)*, vol. 1, no. 1, 2016.
- [4] S. Badawi, "Transformer-Based Neural Network Machine Translation Model for the Kurdish Sorani Dialect," UHD Journal of Science and Technology, vol. 7, no. 1, pp. 15–21, Jan. 2023, doi: 10.21928/uhdjst.v7n1y2023.pp15-21.
- [5] A. Vaswani *et al.*, "Attention Is All You Need," in *Advances in Neural Information Processing System 30 (NIPS 2017)*, Callifornia, United States: Curran Associates, Inc., 2017.
- [6] M. M. Shabrina, A. Aranta, and B. Irmawati, "Rancang Bangun Algoritma Konversi Suara Berbahasa Indonesia Menjadi Teks Latin Berbahasa Sasak Menggunakan Metode Dictionary Based," *Jurnal Teknologi Informasi, Komputer dan Aplikasinya (JTIKA)*, vol. 6, no. 1, pp. 364–375, Mar. 2024, [Online]. Available: http://jtika.if.unram.ac.id/index.php/JTIKA/
- [7] T. Tian, C. Song, J. Ting, and H. Huang, "A French-to-English Machine Translation Model Using Transformer Network," in *Proceedia Computer Science*, Chengdu: Elsevier B.V., 2021, pp. 1438–1443. doi: 10.1016/j.procs.2022.01.182.
- [8] M. Gupta, "Robust Neural Language Translation Model Formulation using Seq2seq approach," *Practice and Applications (FPA)*, vol. 5, no. 2, pp. 61–67, 2021, doi: 10.5281/zenodo.5270391.

- [9] M. Bengi, "Analisis Pengembangan Model Neural Machine Translation (NMT) dengan Transformer untuk Penerjemah Bahasa Indonesia ke Bahasa Gayo," Medan, 2024. [Online]. Available: https://repositori.usu.ac.id/handle/123456789/962 64
- [10] Y. Fauziyah *et al.*, "Mesin Penterjemah Bahasa Indonesia-Bahasa Sunda Menggunakan Recurrent Neural Networks," *JURNAL TEKNOINFO*, vol. 16, no. 2, pp. 313–322, Jul. 2022, [Online]. Available: https://ejurnal.teknokrat.ac.id/index.php/teknoinf
- o/index
 [11] A. Santoso and G. Ariyanto, "Implementasi Deep Learning Berbasis Keras Untuk Pengenalan Wajah," *Jurnal Teknik Elektro*, vol. 18, no. 1, pp. 15–21, Jul. 2018, [Online]. Available: https://www.mathworks.com/discovery/convol
- [12] R. Darma Nurfita and G. Ariyanto, "Implementasi Deep Learning Berbasis Tensorflow Untuk Pengenalan Sidik Jari," *Jurnal Teknik Elektro*, vol. 18, no. 1, pp. 22–27, 2018, [Online]. Available: http://bias.csr.unibo.it/fvc2004/databases.asp
- [13] S. P. Singh, A. Kumar, H. Darbari, L. Singh, A. Rastogi, and S. Jain, "2017 International Conference on Computer, Communications and Electronics (Comptelix)," in 2017 International Conference on Computer, Communication and Electronics (Comptetix), Jaipur: IEEE, Jul. 2017, pp. 162–167.
- [14] R. Juanda and I. Z. Yadi, "Penerapan Rule Based Dengan Algoritma Viterbi Untuk Deteksi Kesalahan Huruf Kapital Pada Karya Ilmiah," *Journal of Computer and Information Systems Ampera*, vol. 1, no. 1, pp. 2775–2496, Jan. 2020, [Online]. Available: https://journalcomputing.org/index.php/journal-cisa/index
- P. Adi Nugroho, I. Fenriana, and R. Arijanto, "Implementasi Deep Learning Menggunakan Convolutional Neural Network (CNN) Pada Ekspresi Manusia," JURNAL ALGOR, vol. 2, no. 1, 2020, [Online]. Available: https://jurnal.buddhidharma.ac.id/index.php/algor /index
- [16] A. Bahari and K. E. Dewi, "Peringkasan Teks Otomatis Abstraktif Menggunakan Transformer Pada Teks Bahasa Indonesia," vol. 13, no. 1, 2024.