

# **Language translator with sequence model**

**A PROJECT REPORT  
for  
AI Project (AI101B)  
Session (2024-25)**

**Submitted by**

**Sandeep kumar  
(202410116100180)  
Preet Kumar  
(202410116100147)**

**Submitted in partial fulfilment of the  
Requirements for the Degree of**

**MASTER OF COMPUTER APPLICATION**

**Under the Supervision of  
Mrs. Komal Salgotra  
Assistant Professor**



**Submitted to**

**DEPARTMENT OF COMPUTER APPLICATIONS  
KIET Group of Institutions, Ghaziabad  
Uttar Pradesh-201206**

# TABLE OF CONTENTS

1. INTRODUCTION.....	2
2. METHODOLOGY.....	3
1. Approach	
2. Algorithm Used	
3. CODE.....	6
4. OUTPUTS SCREENSHOTS .....	11
5. OUTPUT EXPLANATION.....	13

## Introduction

Language translation stands as one of the most impactful and challenging tasks within the domain of Natural Language Processing (NLP). With globalization and the increasing need for multilingual communication, automatic translation systems have become vital tools in breaking language barriers. This project aims to build an English-to-French translation model using a Sequence-to-Sequence (Seq2Seq) architecture powered by Long Short-Term Memory (LSTM) networks — a specialized form of recurrent neural networks (RNNs) well-suited for learning from sequential data.

The task of machine translation involves converting text or speech from one language (source) to another (target) while preserving its meaning, tone, and grammatical correctness. Unlike rule-based or statistical methods, neural machine translation (NMT) leverages deep learning to capture the contextual and semantic intricacies of human language, thereby producing more natural and accurate translations. Among the various deep learning models, the Seq2Seq model has emerged as a powerful framework for handling variable-length input and output sequences, making it ideal for translation tasks.

In this project, the Seq2Seq architecture comprises two main components — an Encoder and a Decoder, both built using LSTM layers. The encoder processes the input English sentence and condenses it into a fixed-size context vector that encapsulates the sentence's meaning. This context vector is then passed to the decoder, which generates the corresponding French sentence, one word at a time. The LSTM units are particularly effective in this setting due to their ability to retain long-term dependencies, manage vanishing gradients, and model the intricate structure of human language over time.

The dataset used for training the model consists of thousands of paired English and French sentences, allowing the model to learn direct mappings between the two languages. The model is trained to minimize the difference between its predicted output and the actual French sentence, using categorical cross-entropy as the loss function and techniques such as teacher forcing to improve learning efficiency.

Through this project, we explore how deep learning can be applied to solve real-world language translation problems. The implementation demonstrates the power of sequence modeling in NLP and lays the foundation for building more advanced models, such as those incorporating attention mechanisms or Transformer-based architectures.

This English-to-French translator serves not only as a practical application of neural networks in translation but also as a comprehensive introduction to the workings of Seq2Seq LSTM models in handling sequential, multilingual data in NLP.

# Methodology

## A. Approach

The project is structured into the following steps to build an English-to-French language translator using a Sequence-to-Sequence model with LSTM networks:

### 1. Data Collection:

Parallel corpora consisting of English sentences and their corresponding French translations are collected from publicly available datasets (e.g., Tatoeba, ManyThings.org, or other bilingual sources). These datasets provide clean, aligned sentence pairs for supervised learning.

### 2. Data Preprocessing:

- **Text Cleaning:** Removing punctuation, converting text to lowercase, and standardizing special characters.
- **Tokenization:** Breaking sentences into individual words or subwords using tokenizer tools.
- **Padding and Sequence Formatting:** Padding sequences to a fixed length and adding special tokens like <start> and <end> to define sentence boundaries.
- **Vocabulary Indexing:** Creating word-to-index and index-to-word dictionaries for both source (English) and target (French) languages.

### 3. Exploratory Data Analysis (EDA):

- **Sentence Length Distribution:** Analyzing the average and maximum sentence lengths in both languages.
- **Vocabulary Size:** Measuring the number of unique tokens in English and French.

- **Word Frequency:** Identifying the most frequently occurring words to assist in vocabulary pruning or optimization.

#### 4. Model Selection & Training:

- **Encoder-Decoder Architecture:** Implementing a Seq2Seq model with two LSTM networks — one for the encoder (English input) and one for the decoder (French output).
- **Teacher Forcing:** During training, the actual target output (French sentence) is fed into the decoder to improve convergence and accuracy.
- **Embedding Layers:** Mapping word indices to dense vectors for better semantic understanding.
- **Training the Model:** Using categorical cross-entropy loss and optimization with Adam or RMSprop optimizers.

#### 5. Evaluation & Translation:

- **Performance Metrics:** Evaluating translation quality using BLEU (Bilingual Evaluation Understudy) scores and perplexity.
- **Sample Translations:** Generating and comparing predicted French sentences with ground truth outputs.
- **Visualization:** Displaying attention weights (if attention mechanism is used) and example translations for interpretability.

---

## B. Algorithms Used

### 1. Sequence-to-Sequence (Seq2Seq) with LSTM

- **Overview:**

Seq2Seq is a deep learning model that handles variable-length input and output sequences. It is ideal for machine translation, where the number of words in source and target sentences can differ.
- **Structure:**
  - **Encoder:** Processes the English sentence and compresses it into a fixed-size context vector (hidden state).
  - **Decoder:** Takes the context vector and generates the French sentence word by word.
- **Key Features:**
  - Uses **LSTM units** to effectively capture long-term dependencies and avoid vanishing gradient issues.
  - Capable of handling different sentence lengths.
  - Optionally extended with an **Attention Mechanism** to improve translation of longer or complex sentences.

## 2. Attention Mechanism (Optional Enhancement)

- **Purpose:**

Allows the decoder to focus on specific parts of the input sentence at each decoding step, improving translation accuracy, especially for longer sentences.
- **Working:**
  - Computes a weighted sum of all encoder hidden states.
  - Aligns each output word with relevant input words dynamically.

# Code

```
import numpy as np
import tensorflow as tf
from tensorflow.keras.models import Model
from tensorflow.keras.layers import Input, LSTM, Dense, Embedding
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences

# 1. Sample dataset
english_texts = [
    'hello', 'how are you', 'good morning', 'good night', 'thank you',
    'please', 'yes', 'no', 'what is your name', 'my name is John',
    'where are you from', 'i am from india', 'do you speak english',
    'i love you', 'see you later', 'i am fine', 'welcome',
    'goodbye', 'excuse me', 'i don't understand'
]

french_texts = [
    'salut', 'comment ça va', 'bonjour', 'bonne nuit', 'merci',
    's'il vous plaît', 'oui', 'non', 'comment tu t'appelles', 'je m'appelle John',
    'd'où viens-tu', 'je viens d'Inde', 'parles-tu anglais',
    'je t'aime', 'à plus tard', 'je vais bien', 'bienvenue',
    'au revoir', 'excusez-moi', 'je ne comprends pas'
]

# Add special tokens
french_texts = ['<start> ' + text + ' <end>' for text in french_texts]

# 2. Tokenization
eng_tokenizer = Tokenizer()
```

```

fre_tokenizer = Tokenizer(filters='') # Preserve <start>, <end>

eng_tokenizer.fit_on_texts(english_texts)
fre_tokenizer.fit_on_texts(french_texts)

eng_sequences = eng_tokenizer.texts_to_sequences(english_texts)
fre_sequences = fre_tokenizer.texts_to_sequences(french_texts)

max_encoder_seq_length = max(len(seq) for seq in eng_sequences)
max_decoder_seq_length = max(len(seq) for seq in fre_sequences)

encoder_input_data = pad_sequences(eng_sequences, maxlen=max_encoder_seq_length,
padding='post')
decoder_input_data = pad_sequences([seq[:-1] for seq in fre_sequences],
maxlen=max_decoder_seq_length - 1, padding='post')
decoder_target_data = pad_sequences([seq[1:] for seq in fre_sequences],
maxlen=max_decoder_seq_length - 1, padding='post')

# 3. Vocabulary sizes
num_encoder_tokens = len(eng_tokenizer.word_index) + 1
num_decoder_tokens = len(fre_tokenizer.word_index) + 1

# 4. Model Parameters
embedding_dim = 64
latent_dim = 256

# 5. Encoder Model
encoder_inputs = Input(shape=(None,))
enc_emb = Embedding(num_encoder_tokens, embedding_dim)(encoder_inputs)
encoder_outputs, state_h, state_c = LSTM(latent_dim, return_state=True)(enc_emb)
encoder_states = [state_h, state_c]

```

```

# 6. Decoder Model

decoder_inputs = Input(shape=(None,))

dec_emb_layer = Embedding(num_decoder_tokens, embedding_dim)

dec_emb = dec_emb_layer(decoder_inputs)

decoder_lstm = LSTM(latent_dim, return_sequences=True, return_state=True)

decoder_outputs, _, _ = decoder_lstm(dec_emb, initial_state=encoder_states)

decoder_dense = Dense(num_decoder_tokens, activation='softmax')

decoder_outputs = decoder_dense(decoder_outputs)

# 7. Define Full Model

model = Model([encoder_inputs, decoder_inputs], decoder_outputs)

model.compile(optimizer='rmsprop', loss='sparse_categorical_crossentropy',
metrics=['accuracy'])

# Expand target data for sparse loss

decoder_target_data = np.expand_dims(decoder_target_data, -1)

# 8. Train the Model

model.fit([encoder_input_data, decoder_input_data], decoder_target_data,
batch_size=16, epochs=300, verbose=1)

# 9. Inference Encoder Model

encoder_model = Model(encoder_inputs, encoder_states)

# 10. Inference Decoder Model

decoder_state_input_h = Input(shape=(latent_dim,))

decoder_state_input_c = Input(shape=(latent_dim,))

decoder_states_inputs = [decoder_state_input_h, decoder_state_input_c]

dec_emb2 = dec_emb_layer(decoder_inputs)

```

```

decoder_outputs2, state_h2, state_c2 = decoder_lstm(dec_emb2,
initial_state=decoder_states_inputs)
decoder_outputs2 = decoder_dense(decoder_outputs2)
decoder_states2 = [state_h2, state_c2]

decoder_model = Model(
    [decoder_inputs] + decoder_states_inputs,
    [decoder_outputs2] + decoder_states2
)

# 11. Reverse Token Dictionaries
reverse_fre_index = dict((i, word) for word, i in fre_tokenizer.word_index.items())
fre_word_index = fre_tokenizer.word_index

# 12. Decode Sequence Function
def decode_sequence(input_seq):
    states_value = encoder_model.predict(input_seq)

    target_seq = np.array([[fre_word_index['<start>']]])
    stop_condition = False
    decoded_sentence = ""

    while not stop_condition:
        output_tokens, h, c = decoder_model.predict([target_seq] + states_value)

        sampled_token_index = np.argmax(output_tokens[0, -1, :])
        sampled_word = reverse_fre_index.get(sampled_token_index, "")

        if sampled_word == '<end>' or len(decoded_sentence.split()) > max_decoder_seq_length:
            stop_condition = True
        else:

```

```

decoded_sentence += ' ' + sampled_word

target_seq = np.array([[sampled_token_index]])
states_value = [h, c]

return decoded_sentence.strip()

# 13. Translate Function

def translate(input_text):

    input_seq = eng_tokenizer.texts_to_sequences([input_text])
    input_seq = pad_sequences(input_seq, maxlen=max_encoder_seq_length, padding='post')
    return decode_sequence(input_seq)

# 14. User Interface

while True:

    input_text = input("\nEnter an English sentence to translate to French (or type 'exit'): ")
    if input_text.lower() == 'exit':
        print("Exiting translator...")
        break
    translation = translate(input_text)
    print("French Translation:", translation)

```

## Outputs

```
Enter an English sentence to translate to French (or type 'exit'): hello
```

```
1/1 ━━━━━━ 0s 193ms/step
```

```
1/1 ━━━━━━ 0s 241ms/step
```

```
1/1 ━━━━━━ 0s 41ms/step
```

```
French Translation: salut
```

```
Enter an English sentence to translate to French (or type 'exit'): welcome
```

```
1/1 ━━━━━━ 0s 38ms/step
```

```
1/1 ━━━━━━ 0s 46ms/step
```

```
1/1 ━━━━━━ 0s 43ms/step
```

```
French Translation: bienvenue
```

```
Enter an English sentence to translate to French (or type 'exit'): i am from india
```

```
1/1 ━━━━━━ 0s 56ms/step
```

```
1/1 ━━━━━━ 0s 56ms/step
```

```
1/1 ━━━━━━ 0s 69ms/step
```

```
1/1 ━━━━━━ 0s 55ms/step
```

```
1/1 ━━━━━━ 0s 67ms/step
```

```
French Translation: je viens d'inde
```

```
Enter an English sentence to translate to French (or type 'exit'): exit
```

```
Exiting translator...
```

```
Epoch 1/300
2/2 5s 49ms/step - accuracy: 0.0642 - loss: 3.6846
Epoch 2/300
2/2 0s 53ms/step - accuracy: 0.3742 - loss: 3.6105
Epoch 3/300
2/2 0s 49ms/step - accuracy: 0.3742 - loss: 3.5052
Epoch 4/300
2/2 0s 48ms/step - accuracy: 0.3825 - loss: 3.2623
Epoch 5/300
2/2 0s 54ms/step - accuracy: 0.3700 - loss: 2.7097
Epoch 6/300
2/2 0s 48ms/step - accuracy: 0.3742 - loss: 2.2813
Epoch 7/300
2/2 0s 51ms/step - accuracy: 0.3825 - loss: 2.1508
Epoch 8/300
2/2 0s 46ms/step - accuracy: 0.3742 - loss: 2.1621
Epoch 9/300
2/2 0s 47ms/step - accuracy: 0.3867 - loss: 2.0445
Epoch 10/300
2/2 0s 48ms/step - accuracy: 0.4500 - loss: 2.1151
Epoch 11/300
2/2 0s 50ms/step - accuracy: 0.3792 - loss: 2.0831
Epoch 12/300
2/2 0s 50ms/step - accuracy: 0.3925 - loss: 2.1067
Epoch 13/300
2/2 0s 66ms/step - accuracy: 0.4000 - loss: 1.9730
Epoch 14/300
2/2 0s 48ms/step - accuracy: 0.5692 - loss: 1.9675
Epoch 15/300
2/2 0s 53ms/step - accuracy: 0.5700 - loss: 2.0112
Epoch 16/300
2/2 0s 50ms/step - accuracy: 0.4742 - loss: 1.9282
Epoch 17/300
2/2 0s 49ms/step - accuracy: 0.5800 - loss: 1.8608
Epoch 18/300
2/2 0s 48ms/step - accuracy: 0.5758 - loss: 1.8611
Epoch 19/300
2/2 0s 49ms/step - accuracy: 0.5758 - loss: 1.8417
Epoch 20/300
2/2 0s 51ms/step - accuracy: 0.5758 - loss: 1.8418
Epoch 21/300
2/2 0s 51ms/step - accuracy: 0.5783 - loss: 1.8825
Epoch 22/300
2/2 0s 48ms/step - accuracy: 0.5825 - loss: 1.8047
Epoch 23/300
2/2 0s 49ms/step - accuracy: 0.5867 - loss: 1.8066
Epoch 24/300
2/2 0s 51ms/step - accuracy: 0.5783 - loss: 1.8185
Epoch 25/300
2/2 0s 47ms/step - accuracy: 0.5742 - loss: 1.7946
Epoch 26/300
2/2 0s 49ms/step - accuracy: 0.5825 - loss: 1.7525
Epoch 27/300
2/2 0s 47ms/step - accuracy: 0.5825 - loss: 1.7323
Epoch 28/300
2/2 0s 50ms/step - accuracy: 0.5700 - loss: 1.7687
```

## OUTPUT EXPLANATION

### Screenshot 1: Translation in Action

This screenshot represents the inference (or prediction) phase of your Seq2Seq (Sequence-to-Sequence) LSTM (Long Short-Term Memory) model, which translates English sentences into French.

Let's break it down:

- Input/Output Examples: The table shows several translation examples, each with:
  - English Input: The sentence in English.
  - French Output: The corresponding translation in French.
  - Status: Whether the translation is correct or if there's an issue.

For example:

- "hello" → "salut":  This is a correct translation.
- "welcome" → "bienvenue":  Correct translation.
- "i am from india" → "je viens d'inde":  The translation is correct.
- "exit" → "Exits the program":  This shows that the model correctly handles special cases like a command (i.e., exiting the program).

These results indicate that the model is successfully translating the input sentences, but keep in mind, this performance is based on the small dataset it was trained on.

- Time per Step: The line "1/1 ===== 0s 56ms/step" shows the inference speed. Since you're only translating one sentence at a time, each step of the translation process is very quick. The model is making a prediction for a single input (1/1), and the time taken per step is 56 milliseconds.

### Screenshot 2: Model Training Progress

This screenshot provides an overview of the model's training performance over 300 epochs. It gives us insights into how well the model is learning over time.

- Loss Decrease:
  - Loss is a measure of how well the model's predictions match the target (in this case, the French translation). A lower loss indicates that the model is getting better at making correct predictions.
  - The loss starts at 3.68 in epoch 1, which is relatively high, meaning the model isn't performing well initially. However, by epoch 28, the loss decreases to around 1.76, showing improvement in the model's ability to generate accurate translations.

In machine learning, the goal is to minimize loss, and the steady decrease here indicates that the model is learning over time.