Project Documentation: Arabic-English Translation using MarianMT

Project Title

Arabic-English Neural Machine Translation using MarianMT

Team Member Names:

- Yassmin ahmed ibrahim
- Maram essam ahmed
- Youssef tarek aboelfotouh
- Warda Khaled dahy
- Mariam ahmed Mostafa

Youssef Mohamed fathy

Overview

This project focuses on building a neural machine translation (NMT) system that translates English text into Arabic using the pre-trained MarianMT model from Hugging Face. The system is trained and evaluated using a parallel corpus of English-Arabic sentence pairs. The implementation is based on the transformers library, and the training pipeline leverages Hugging Face's Seq2SeqTrainer.

1. Dataset Information

Dataset Source: Kaggle

• Name: Arabic to English Translation Sentences

- . File Used: ara eng.txt
- . Columns: english, arabic
- Format: Tab-separated values (TSV)

Dataset Preprocessing

- Removed punctuation and unnecessary characters.
- Lowercased English text.
- Cleaned Arabic text using Unicode ranges for Arabic.
- Filtered out sentence pairs with less than 3 words or less than 5 characters.
- Removed duplicate English entries.

Data Split

. Training Set: 80%

• Testing Set: 20%

Preprocessing Steps

.English Text:

Converted text to lowercase.

Removed brackets and their content: (...), [...].

Removed unnecessary characters while keeping useful punctuation: ., ,, !, ?.

Normalized characters using unicodedata.normalize() to standardize accented letters.

.Arabic Text:

Removed all non-Arabic characters using Unicode range: \u0600-\u06FF.

Removed diacritics and decorative marks, such as: -.

Removed digits and non-alphabetic symbols.

2. Model Information

Pre-trained Model

- . Model Name: Helsinki-NLP/opus-mt-tc-big-en-ar
- . Library: transformers by Hugging Face

• Architecture: MarianMT (based on Transformer architecture)

Tokenizer

. Type: SentencePiece

. Tokenizer: MarianTokenizer from

Hugging Face

Input Length

• Maximum Sequence Length: 128 tokens

3. Training Details

Preprocessing

- Tokenization applied to both English input and Arabic target.
- Padding and truncation used for fixed-length inputs (max length = 128).

Training Configuration

. Epochs: 10

. Training Batch Size: 16

. Evaluation Batch Size: 16

• Output Directory: ./results

. Logging Directory: ./logs

- Logging Steps: 10
- **Save Steps:** 1000
- Evaluation Strategy: Custom bertscore callback
- Prediction: Enabled with
 predict with generate=True

Training Loss (last few steps):

- Decreased from ~ 6.5 to ~ 0.2 over 10 epochs
- Indicates effective learning and convergence

4. Evaluation Metrics

Metric Used: BERTScore

- **Library:** bert_score (from HuggingFace)
- Language Models Used: Pre-trained multilingual BERT models (bert-base-multilingual-cased)
- Method:
 - Computes similarity between reference and generated translation using contextual embeddings
 - Captures semantic meaning better than lexical matchbased metrics like BLEU
- Evaluation Frequency:
 - Evaluated on 100 random test samples at the end of each epoch
 - o Provides F1 score for precision-recall balance
- Why BERTScore?
 - o More reliable for low-resource language pairs
 - Better reflects human judgment of translation quality

5. Model Limitations

. Limited Domain Generalization:

 Trained on a small parallel corpus; may not generalize to other domains (e.g., medical, legal).

. Token Truncation:

 Long sentences may be truncated, potentially affecting translation accuracy.

· Vocabulary Limitation:

 Pre-trained tokenizer might not handle rare or domain-specific tokens well.

Cultural and Contextual Nuances:

 May miss idiomatic expressions or cultural references that require contextual understanding.

• BERTScore Reliability:

- While more semantically aware than BLEU, still limited by the capabilities of the underlying BERT model
- Sensitive to language-specific pretraining and sentence structure nuances.

6.Parameters Model

1. model_name

o **Type:** str

 Description: The name of the pre-trained Marian model used for English-to-Arabic translation. It is used to load the model and tokenizer.

2. max length

- o Type: int
- Description: The maximum token length for the input and output sequences during tokenization.

train_dataset & test_dataset

- o Type: Dataset
- Description: The training and testing datasets that are used for training and evaluation. These are instances of the Dataset class from the datasets library.

4. training_args

- o **Type:** Seq2SeqTrainingArguments
- Description: This contains all the hyperparameters required for training the model. It specifies batch size, number of epochs, logging settings, save steps, etc.

o Important fields:

- per_device_train_batch_size: The batch size for training.
- num_train_epochs: The number of epochs for training.
- logging_dir: Directory where logs will be saved.
- save_steps: Frequency at which the model is saved.

5. generation_config

- o **Type:** GenerationConfig
- Description: Configuration for controlling the generation of text by the model.

o Important fields:

- max_length: Maximum length of the generated sequence.
- num_beams: Number of beams for beam search (controls the search size during translation).
- length_penalty: Penalty for longer translations.
- early_stopping: Whether to stop generating early if the sequence is finished.

 bad_words_ids: A list of words to avoid in the generated text.

6. batch_size

- o Type: int
- Description: The number of samples in each batch during translation.

7. inputs (in preprocess function)

- o Type: dict
- Description: This is the tokenized input for the English text. It is a dictionary containing keys like input_ids and attention_mask, which are used by the model to process the input.

8. labels (in preprocess function)

- o **Type:** dict
- Description: The tokenized version of the Arabic text (target). It also contains input_ids which represent the tokens for the target language.

9. decoded preds and decoded labels

- o **Type:** list of str
- o **Description:** The translated output (decoded_preds) and the reference translation (decoded_labels). These are lists of strings containing the decoded translations.

10. P, R, F1

- o **Type:** Tensor
- Description: The Precision (P), Recall (R), and F1 scores computed using the BERTScore metric to evaluate the translation quality.

11. preds and labels (in compute_metrics)

- o **Type:** list of int
- Description: These are the predicted and true label sequences (after tokenization) used for computing metrics like BERTScore.

12. texts (in generate_translation)

- o **Type:** list of str
- Description: The English text to be translated. This is passed in batches to the model for translation.
- 13. input_ids and attention_mask (in generate_translation)

- o **Type:** Tensor
- o **Description:** These are the inputs to the model during generation. input_ids are the tokenized representation of the input texts, and attention_mask is used to indicate which tokens should be attended to (ignores padding tokens).

14. interface (in Gradio UI)

- o **Type:** gr.Interface
- Description: This is the Gradio interface that takes an English sentence as input and displays the predicted Arabic translation as output. It uses the generate_translation function as its backend processing logic.

7. Future Enhancements

Data Augmentation:

 Use back-translation or synonym replacement to enrich dataset.

. Fine-Tuning:

 Use a larger or more domain-specific dataset.

. UI Integration:

 Integrate with a GUI (e.g., Gradio or Streamlit) for easy testing and deployment.

. Deploy API:

 Wrap the model in a Flask or FastAPI service and deploy to Vercel or Hugging Face Spaces.

8. Dependencies

- transformers
- datasets
- sentencepiece
- scikit-learn
- nltk
- torch
- bert_score

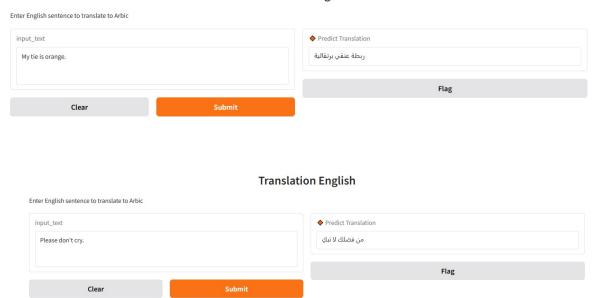
9. Potential Improvements

- Incorporate domain-specific data for fine-tuning
- Use larger context-aware models (e.g., mBART or T5 multilingual)
- Introduce post-processing grammar correction or reranking
- Utilize semantic-aware metrics (e.g., BERTScore) alongside BLEU

10.RUN

Translation English Enter English sentence to translate to Arbic input_text I have to go home. Clear Submit

Translation English



11. Conclusion

This project successfully demonstrates the application of a pre-trained MarianMT model for Arabic-English translation. Through fine-tuning and bertscore evaluation, the model achieved good performance on a general dataset, with opportunities for further enhancement in specialized applications.

•