



TRANSLATION

TEAM MEMBERS

Abdalrahman Abdelaziz Ebrahim Ebrahim	20210514
Mohamed Hassan Mohamed Hafez	20210762
Abdullah Mahmoud Abdulmohsen Ahmed	20210560
Hager Sayed Rashid	20211023
Huda Maher Mohamed Mohamed	20211031
Mona Mohamed Abdelaziz	20210967

PROJECT OVERVIEW

This project implements a bidirectional Arabic-English neural machine translation system using transformer-based models from the MarianMT framework. The system provides high-quality translations in both directions while being computationally efficient.

The background is a dark blue gradient. On the left and right sides, there are decorative elements: a series of overlapping, semi-transparent blue squares and rectangles, and a pattern of small, light blue dots arranged in a grid. The text "TECHNICAL ARCHITECTURE" is centered in the middle of the image in a bold, white, sans-serif font.

TECHNICAL ARCHITECTURE

MODEL SPECIFICATIONS

- **BASE MODEL: HELSINKI-NLP MARIANMT (TRANSFORMER ARCHITECTURE)**

- **ARABIC-TO-ENGLISH MODEL:**

LAYERS: 6 ENCODER, 6 DECODER

HIDDEN SIZE: 512

ATTENTION HEADS: 8

PARAMETERS: ~77 MILLION

FINETUNED ON CUSTOM DATASET

- **ENGLISH-TO-ARABIC MODEL:**

PRETRAINED HELSINKI-NLP MODEL

SAME ARCHITECTURE AS ABOVE

USED WITHOUT FINE-TUNING

TRAINING DETAILS

- **FRAMEWORK: HUGGINGFACE TRANSFORMERS**

- **TRAINING ARGS:**

BATCH SIZE: 8 (PER DEVICE)

LEARNING RATE: 5E-5

EPOCHS: 3

MIXED PRECISION (FP16)

- **HARDWARE: GPU-ACCELERATED (GOOGLE COLAB)**

DATASET INFORMATION

- **SOURCE DATA**

BILINGUAL SENTENCE PAIRS FROM OPEN-SOURCE PARALLEL CORPUS

CONTAINS 10,000 ARABIC-ENGLISH PAIRS

BALANCED ACROSS DOMAINS (NEWS, CONVERSATIONS, GENERAL TEXT)

- **LINK: [NMT-WITH-ATTENTION-FOR-AR-TO-EN](#)**

DATASET INFORMATION

- **PREPROCESSING PIPELINE**

UNICODE NORMALIZATION

CASE NORMALIZATION

SPECIAL TOKEN HANDLING

PUNCTUATION STANDARDIZATION

TEXT CLEANING REGEX PATTERNS

SEQUENCE BOUNDARY MARKING (<START>, <END>)

- **DATA SPLITS**

TRAINING: 90% (8,000 PAIRS)

TEST: 10% (2,000 PAIRS)

The background is a dark blue gradient. On the left and right sides, there are decorative elements: a series of overlapping, semi-transparent blue squares that create a sense of depth and movement, and a pattern of small, light blue dots arranged in a grid. The text "PERFORMANCE METRICS" is centered in the middle of the image in a bold, white, sans-serif font.

PERFORMANCE METRICS

QUANTITATIVE EVALUATION



- BLEU SCORE: 0.72 (ARABIC→ENGLISH)

0-1 SCALE WHERE HIGHER IS BETTER

COMPETITIVE WITH BASELINE MODELS

HUMAN TRANSLATION TYPICALLY SCORES 0.6-0.7



QUALITATIVE ASSESSMENT

- HANDLES COMMON PHRASES WELL
- MAINTAINS REASONABLE GRAMMAR
- PRESERVES MOST SEMANTIC MEANING

SYSTEM LIMITATIONS

- **TECHNICAL CONSTRAINTS:**

MAX SEQUENCE LENGTH: 128 TOKENS

LIMITED DOMAIN ADAPTATION

SUBOPTIMAL HANDLING OF:

PROPER NOUNS

IDIOMATIC EXPRESSIONS

HIGHLY TECHNICAL TERMS

- **LINGUISTIC CHALLENGES:**

ARABIC MORPHOLOGICAL COMPLEXITY

DIALECTAL VARIATIONS NOT HANDLED

GENDER/CASE AGREEMENT ISSUES

FORMALITY LEVELS NOT PRESERVED

CONCLUSION

In this project, we developed a bidirectional Arabic-English neural machine translation system using MarianMT transformer models from Hugging Face. Through fine-tuning on a curated parallel corpus, we achieved a competitive BLEU score of 0.72 for Arabic-to-English translation while leveraging a pretrained model for English-to-Arabic translation.

This project highlights the potential of open-source NMT models for bridging language barriers. While challenges remain in handling nuances like idiomatic expressions and dialects, the system provides a strong foundation for future improvements. By refining the model and expanding its training data, we can move closer to human-like translation quality in real-world applications.

REFERENCES

- MarianNMT: <https://marian-nmt.github.io/>
- Hugging Face Transformers: <https://huggingface.co/docs/transformers/index>
- Helsinki-NLP Models: <https://huggingface.co/Helsinki-NLP>
- Dataset Link: [nmt-with-attention-for-ar-to-en](#)

The background is a dark navy blue. On the left side, there are several overlapping, semi-transparent blue geometric shapes, including rectangles and parallelograms, some of which are tilted. On the right side, there is a grid of small, light blue dots that fades out towards the center.

THANK YOU