

# Machine Translation with Marian Transformers: English-to-French Translation Pipeline Analysis

Junaid Ali

June 22, 2025

## Abstract

This report analyzes a machine translation implementation using the Marian transformer architecture for English-to-French translation. The analysis covers dataset preprocessing, model configuration, training methodology, and evaluation metrics, highlighting both strengths and areas for improvement in the implementation.

## 1 Introduction

Machine translation represents one of the most challenging tasks in natural language processing, requiring models to understand semantic meaning across languages while maintaining grammatical correctness and contextual relevance. This analysis examines a practical implementation using the Helsinki-NLP Marian model, a state-of-the-art transformer architecture specifically designed for translation tasks.

## 2 Technical Architecture

### 2.1 Model Selection

The implementation utilizes the `Helsinki-NLP/opus-mt-en-fr` model, which offers several advantages:

- Pre-trained on large-scale parallel corpora
- Optimized specifically for English-French translation
- Compact architecture suitable for fine-tuning
- Strong baseline performance on standard benchmarks

### 2.2 Dataset Configuration

The training dataset comprises:

Training samples = 200,000	(1)
Validation samples = 100,000	(2)
Train-validation ratio = 80 : 20	(3)

## 3 Preprocessing Pipeline

### 3.1 Tokenization Strategy

The preprocessing function implements comprehensive text handling:

```
1 def preprocess_function(examples):
2     inputs = [str(text) if text is not None else ""
3               for text in examples['en']]
4     targets = [str(text) if text is not None else ""
5                for text in examples['fr']]
6
7     model_inputs = tokenizer(
8         inputs,
9         max_length=64,
10        truncation=True,
11        padding='max_length'
12    )
```

#### Critical Analysis:

- **Strengths:** Robust null value handling, consistent padding strategy
- **Weaknesses:** Fixed sequence length may truncate important content
- **Recommendation:** Implement dynamic padding or analyze sequence length distribution

## 4 Training Methodology

### 4.1 Optimization Configuration

The training employs the following hyperparameters:

Parameter	Value
Learning Rate	$5 \times 10^{-5}$
Batch Size	8
Optimizer	AdamW
Precision	FP16
Training Steps	10

Table 1: Training Configuration Parameters

### 4.2 Loss Function Analysis

The training loss progression demonstrates:

$$\text{Initial Loss} \approx 1.63 \quad (4)$$

$$\text{Final Loss} \approx 0.99 \quad (5)$$

$$\text{Average Loss} = 1.87 \quad (6)$$

## 5 Evaluation Metrics

### 5.1 BLEU Score Assessment

The model achieves a BLEU score of **0.6811** on the test translation, indicating:

- Good semantic preservation
- Reasonable grammatical structure
- Minor issues with idiomatic expressions

### 5.2 Qualitative Analysis

**Sample Translation:**

*Input:* "I didn't expect the big departure..."

*Output:* "Je ne m'attendais pas au grand départ..."

*Quality:* High semantic accuracy with minor grammatical variations

## 6 Critical Assessment

### 6.1 Implementation Strengths

1. Comprehensive environment setup with CUDA support
2. Proper dataset splitting and preprocessing
3. Efficient memory usage with FP16 precision
4. Robust error handling for null values
5. Clear progress tracking and logging

## 7 Conclusion

The implementation demonstrates a solid foundation for machine translation using Marian transformers. While the technical setup is comprehensive and the preprocessing pipeline is robust, the extremely limited training duration prevents meaningful model improvement. The achieved BLEU score of 0.68 indicates reasonable translation quality, but the evaluation methodology requires enhancement for production deployment.

The code serves as an excellent educational example but needs significant modifications for practical applications, particularly in training duration, evaluation comprehensiveness, and data analysis depth.