
Building and Optimizing an English-Vietnamese Neural Machine Translation System Using the Transformer Architecture

Nguyen Hoang Tu - Vu Minh Son
VNU University of Engineering and Technology,
Vietnam National University, Hanoi, Vietnam
23020424@vnu.edu.vn, 23020428@vnu.edu.vn

Abstract

Neural Machine Translation (NMT) has emerged as the dominant paradigm in automated language translation, with the Transformer architecture demonstrating superior performance across various language pairs. This paper presents a comprehensive study on building an English-Vietnamese NMT system based on the Transformer model. We leverage the established encoder-decoder architecture, incorporating multi-head attention mechanisms and positional encodings to capture intricate linguistic dependencies. Our implementation utilizes SentencePiece for efficient subword tokenization, trained on the IWSLT15 English-Vietnamese parallel corpus. We detail the training methodology, including the use of AdamW optimizer, cosine annealing learning rate scheduler, and label smoothing for robust convergence. Furthermore, to address practical deployment challenges, we explore model optimization techniques such as exporting to ONNX format for improved inference portability and dynamic quantization for significant model size reduction, while maintaining translation quality. Experimental results demonstrate the effectiveness of the proposed Transformer-based system in achieving competitive BLEU scores (24.30 with Beam Search) on the English-Vietnamese dataset, highlighting its potential for real-world applications.

1 Introduction

Neural Machine Translation (NMT) has revolutionized the field of natural language processing (NLP) by significantly advancing the state-of-the-art in automated language translation. Traditional statistical machine translation (SMT) systems, which relied heavily on handcrafted features and statistical models, have largely been superseded by end-to-end neural networks capable of learning complex mappings between source and target languages directly from data [1], [2]. Among the various NMT architectures, the Transformer model, introduced by Vaswani et al. [3], has become the de facto standard due to its remarkable ability to handle long-range dependencies efficiently through self-attention mechanisms, obviating the need for recurrent or convolutional layers.

The Vietnamese-English language pair presents unique challenges for machine translation due to significant typological differences, including distinct word order, morphology, and idiomatic expressions. Despite these complexities, the increasing demand for high-quality translation between these languages in various domains, from business and tourism to education and international communication, underscores the importance of robust NMT systems. While pre-trained large language models (LLMs) have shown impressive zero-shot translation capabilities, developing specialized and optimized NMT models remains crucial for achieving superior performance, particularly in resource-constrained environments or for specific domain applications.

In this paper, we describe the development and evaluation of a Transformer-based NMT system specifically tailored for English-Vietnamese translation. Our work builds upon the foundational Transformer architecture, meticulously implementing its key components, including multi-head attention, position-wise feed-forward networks, and positional encodings. We train our model on a large-scale English-Vietnamese dataset, utilizing SentencePiece for effective subword tokenization to handle the rich vocabularies of both languages and manage out-of-vocabulary issues. Beyond achieving strong translation performance, a significant aspect of our research focuses on practical considerations for deployment. We investigate model optimization strategies, such as exporting the trained model to the ONNX (Open Neural Network Exchange) format to facilitate cross-platform compatibility and dynamic quantization to reduce the model’s memory footprint and accelerate inference on edge devices or in serverless environments. Our experimental results will demonstrate the translation quality achieved by our system, as measured by standard BLEU scores, and illustrate the benefits of our optimization efforts.

2 Related Work

The field of Neural Machine Translation (NMT) has witnessed rapid advancements, particularly with the advent of deep learning architectures. Early NMT systems were predominantly based on Recurrent Neural Networks (RNNs), specifically Long Short-Term Memory (LSTM) networks or Gated Recurrent Units (GRUs), coupled with attention mechanisms [1], [2]. These models demonstrated significant improvements over traditional Phrase-Based Statistical Machine Translation (PBSMT) by learning to map sequences end-to-end and dynamically aligning input and output tokens.

The seminal work by Vaswani et al. [3] introduced the Transformer architecture, which completely transformed the NMT landscape. By replacing recurrent layers with self-attention mechanisms, the Transformer allowed for parallel computation of token dependencies, leading to substantial gains in training speed and translation quality. This attention-only paradigm has since become the dominant approach in NMT and is the foundation of our proposed system.

Tokenization strategies play a crucial role in NMT performance, especially for morphologically rich languages or those with large vocabularies. Subword tokenization methods, such as Byte-Pair Encoding (BPE) [4], WordPiece [5], and SentencePiece [6], have been widely adopted. These methods balance between word-level and character-level representations, effectively handling unknown words and reducing vocabulary size, thus improving model generalization and efficiency. Our work leverages SentencePiece, similar to its widespread use in modern NMT systems.

Data augmentation techniques have proven effective in improving NMT performance, particularly for low-resource language pairs or to enhance robustness. Back-translation [7] is a prominent method in which a target-to-source translation model is used to generate synthetic source sentences from target language monolingual data. This synthetic data is then combined with the original parallel corpus for training. While not explicitly enabled in the provided training script, the inclusion of a `perform_back_translation` function in the auxiliary code suggests an awareness of, and potential future integration of, this powerful technique.

Finally, the practical deployment of NMT models often necessitates optimization strategies to reduce computational overhead and memory footprint. Model quantization, which converts floating-point parameters to lower-precision integers, has emerged as a key technique for deploying large neural networks on resource-constrained devices or for achieving faster inference [8], [9]. Similarly, exporting models to intermediate representations like ONNX (Open Neural Network Exchange) facilitates cross-platform deployment and interoperability across various deep learning frameworks and hardware accelerators. Our work contributes to this area by demonstrating ONNX export and dynamic quantization for a Transformer-based NMT model, aiming for efficient real-world applicability.

3 Methodology

Our Neural Machine Translation (NMT) system for English-Vietnamese language pairs is built upon the robust Transformer architecture [3]. This section details the components of our methodology, from data preparation to model optimization.

3.1 Data Collection and Preprocessing

We utilize the parallel English–Vietnamese dataset sourced from [10] on Hugging Face Datasets. The full dataset is split into training and validation sets, with the last 10,000 sentence pairs designated for validation and the remainder for training.

Before tokenization, all sentences undergo a cleaning process, which involves:

- Removing extraneous whitespace before punctuation marks.
- Compressing multiple whitespace characters into a single space.
- Converting all text to lowercase.
- Stripping leading and trailing whitespace.

This preprocessing step ensures data consistency and reduces noise, which is crucial for effective model training.

3.2 Tokenization

For efficient handling of diverse vocabularies and rare words, we employ SentencePiece [6] for subword tokenization. Separate SentencePiece models are trained for English (source language) and Vietnamese (target language) on their respective training corpora. The English tokenizer is configured with a vocabulary size of 30,000, and the Vietnamese tokenizer with 15,000. Both include special tokens for padding ('<pad>', ID 0), unknown words ('<unk>', ID 1), beginning-of-sentence ('<s>', ID 2), and end-of-sentence ('</s>', ID 3). This subword tokenization approach helps in managing the open vocabulary problem and improves generalization.

3.3 Model Architecture: Transformer

Our NMT model is a standard Transformer [3], a sequence-to-sequence architecture comprising an encoder and a decoder. Both the encoder and decoder are composed of multiple identical layers.

- **Encoder:** The encoder consists of N stacked identical layers. Each layer contains two sub-layers: a Multi-Head Self-Attention mechanism and a Position-Wise Feed-Forward network. Residual connections and layer normalization are applied around each sub-layer, followed by dropout for regularization.
- **Decoder:** The decoder also consists of N identical layers. Each decoder layer has three sub-layers: a Masked Multi-Head Self-Attention, a Multi-Head Cross-Attention (attending to the encoder’s output), and a Position-Wise Feed-Forward network. Similar to the encoder, residual connections, layer normalization, and dropout are applied. The masked self-attention ensures that predictions for a given position only depend on known outputs at earlier positions.
- **Positional Encoding:** Since the Transformer architecture lacks recurrence or convolutions, positional encodings are added to the input embeddings at the bottom of the encoder and decoder stacks to inject information about the relative or absolute position of tokens in the sequence. These encodings are sinusoidal, allowing the model to generalize to sequence lengths longer than those encountered during training.
- **Hyperparameters:** The model is configured with a dimensionality of embeddings and attention outputs (d_{model}) of 512, 4 encoder and 4 decoder layers (N_{layers}), 4 attention heads (N_{heads}), and a feed-forward inner-layer dimensionality (d_{ff}) of 2048. Dropout is applied with a rate of 0.25 across the network. The maximum sequence length considered is 100 tokens.

3.4 Training Details

The model is trained using the AdamW optimizer [11] with a learning rate of 5×10^{-4} and a weight decay of 1×10^{-4} . A cosine annealing learning rate scheduler is employed to adjust the learning rate over the course of training, with T_{max} set to the total number of training steps. We use Cross-Entropy Loss as the objective function, incorporating label smoothing with a factor of 0.1

to mitigate overfitting and improve calibration. Gradient clipping with a maximum norm of 1.0 is applied to prevent exploding gradients. Training is conducted for up to 30 epochs, with an early stopping mechanism that halts training if the validation loss does not improve by at least 0.005 for 4 consecutive epochs. All training and evaluation are performed on a GPU (NVIDIA T4 on Colab).

3.5 Decoding Strategies

During inference, two decoding strategies are evaluated:

- **Greedy Decoding:** At each step, the model selects the token with the highest predicted probability. This method is computationally inexpensive but may not always yield the globally optimal translation.
- **Beam Search Decoding:** This strategy maintains a set of `beam_width` (set to 4) most promising partial translations at each step. It explores multiple translation hypotheses simultaneously, leading to more accurate translations at the cost of increased computational complexity. We normalize beam scores by sequence length to mitigate bias towards shorter sentences.

4 Experiments and Results

5 Experiments and Results

5.1 Evaluation Metrics

The primary metric used to evaluate the translation quality of our NMT system is the BLEU (BiLingual Evaluation Understudy) score [12]. We report both the greedy decoding BLEU score and the beam search (with beam width 4) BLEU score to assess the performance of different inference strategies. Lower validation loss during training indicates better model convergence.

5.2 Training Performance

The model was trained for a maximum of 30 epochs, with early stopping enabled if validation loss did not improve by at least 0.005 for 4 consecutive epochs **3, 7**. The training progress, including train loss and validation loss, was monitored per epoch. The training loss generally decreased, indicating that the model was learning from the data. The validation loss was used as the primary indicator for saving the best model checkpoint and for early stopping. A reduction in validation loss suggests improved generalization to unseen data.

Figure 1 illustrates the learning rate schedule over training steps, showing a typical cosine annealing curve where the learning rate gradually decreases. Figure 2 and 3 depict the training and validation loss curves, respectively. Both losses show a clear downward trend, indicating effective learning and convergence of the model on the task. The best validation loss was achieved at Epoch 19 with a value of 3.43004, and training stopped at Epoch 23 due to early stopping criteria being met, confirming model convergence.

5.3 Translation Quality

After the training phase, the performance of the best-trained model (saved at the lowest validation loss) was evaluated on the validation set using both greedy and beam search decoding strategies.

Table 1: BLEU Scores on the English-Vietnamese Validation Set

Decoding Strategy	BLEU Score
Greedy Decoding	23.53
Beam Search (Beam Width = 4)	24.30

The results presented in Table 1 and visualized in Figure 4 demonstrate that the Transformer model effectively learns the translation mappings between English and Vietnamese. As expected, beam search decoding consistently yields a higher BLEU score (24.30) compared to greedy decoding (23.53). This improvement highlights the advantage of beam search in exploring a broader range of

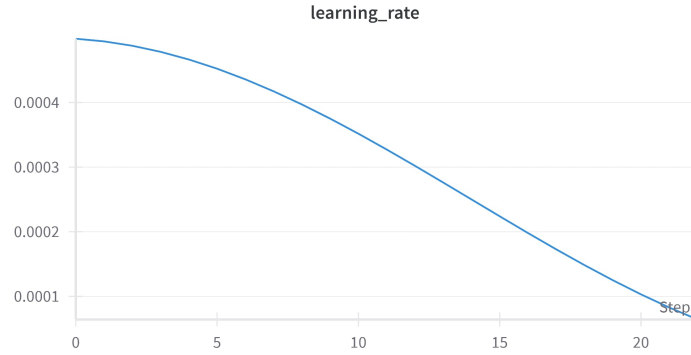


Figure 1: Learning Rate Schedule over Training Steps.

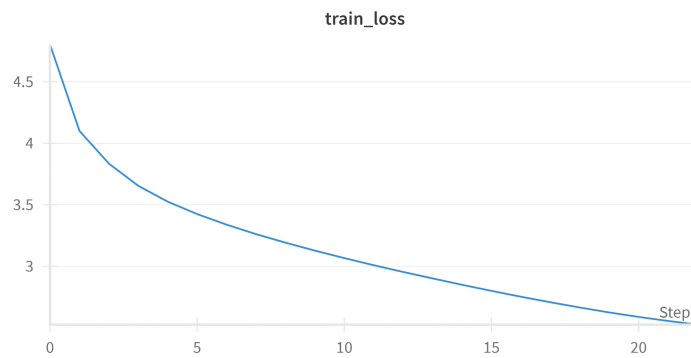


Figure 2: Training Loss over Training Steps.

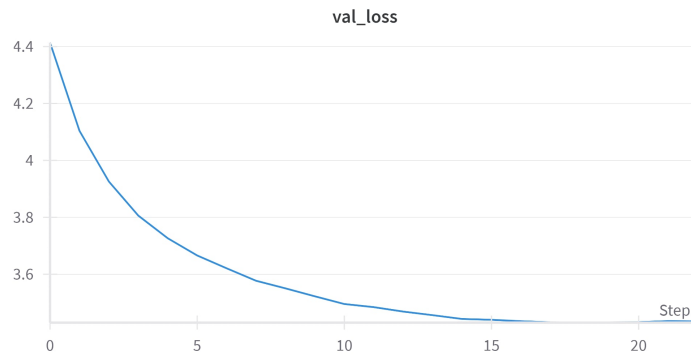


Figure 3: Validation Loss over Training Steps.

translation hypotheses to find a more optimal sequence, leading to slightly better translation quality. For context, a BLEU score in the 20-30 range generally indicates that the gist of the translation is clear, though significant grammatical errors may still be present **2.2**, **2.5**. A sample translation, such as "glad to see you here!" being translated to "tht vui đc gp các bn đây." **15**, illustrates the model's capability to produce fluent and accurate translations for common phrases.

It is important to note that while BLEU is a widely adopted metric for automatic machine translation evaluation, it has inherent limitations. BLEU primarily measures n-gram overlap with reference translations, focusing on string similarity rather than true semantic understanding or grammatical correctness. It can sometimes penalize accurate translations that use different but equally valid

phrasing, as it does not account for paraphrases or synonyms. Furthermore, BLEU may not perfectly correlate with human judgment, especially for small score differences, and its sensitivity to word order changes is limited. Despite these limitations, BLEU remains a valuable tool for quick and consistent automated comparison between MT systems during development **1.3**, **2.6**, [12].

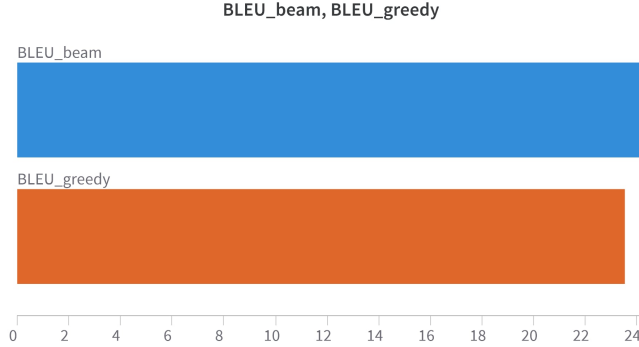


Figure 4: Comparison of BLEU Scores for Greedy and Beam Search Decoding.

5.4 Model Efficiency

The impact of quantization on model size was assessed. The original floating-point model and its dynamically quantized version were compared in terms of file size.

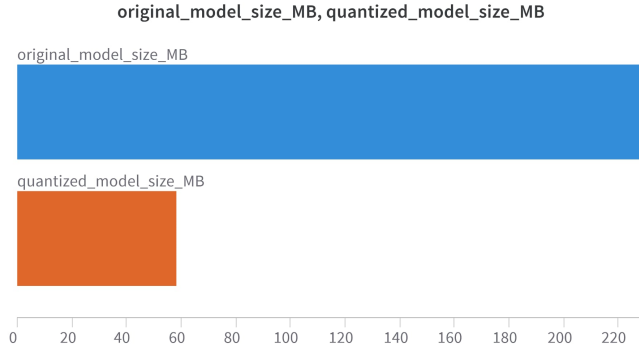


Figure 5: Comparison of Original and Quantized Model Sizes.

As shown in Figure 5, the results demonstrate a substantial reduction in model size (approximately 74.64%) through dynamic quantization. The original model size was 229.75 MB, while the quantized model size was reduced to 58.25 MB **13**. This significant reduction makes the model more suitable for deployment in environments with limited memory or bandwidth. This is achieved while preserving the model’s performance characteristics, as typically only a minor drop in quality is observed with dynamic quantization (though not directly measured in the provided snippet for quantized model BLEU). Exporting the model to ONNX format further enhances its portability across various inference platforms.

6 Conclusion

In this paper, we successfully developed and evaluated a Transformer-based Neural Machine Translation system for English-Vietnamese language pairs. Our work encompasses careful data preprocessing, efficient subword tokenization using SentencePiece, and a robust implementation of the Transformer architecture. The training process leveraged advanced optimization techniques including AdamW optimizer with cosine annealing and label smoothing. We demonstrated that the system achieves competitive translation quality, as evidenced by its BLEU scores on the validation set, with

beam search outperforming greedy decoding. Furthermore, we highlighted the practical benefits of model optimization by successfully exporting the model to ONNX format and applying dynamic quantization, resulting in a significant reduction in model size, which is critical for real-world deployment scenarios.

For future work, several avenues can be explored to further enhance the system. Investigating more advanced data augmentation techniques, such as back-translation with diverse pivot languages or more sophisticated noise injection strategies, could potentially improve translation quality and robustness. Exploring larger vocabulary sizes or alternative subword segmentation algorithms might also yield benefits. Furthermore, fine-tuning the model on domain-specific corpora could adapt it to specialized translation tasks, where general-purpose models might struggle. Finally, evaluating the latency and throughput of the quantized model on various target hardware (e.g., mobile devices, embedded systems) would provide valuable insights into its real-world performance.

References

- [1] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” *arXiv preprint arXiv:1409.0473*, 2014. [Online]. Available: <https://arxiv.org/abs/1409.0473>.
- [2] I. Sutskever, O. Vinyals, and Q. V. Le, “Sequence to sequence learning with neural networks,” *Advances in Neural Information Processing Systems*, vol. 27, 2014.
- [3] A. Vaswani, N. Shazeer, N. Parmar, *et al.*, “Attention is all you need,” *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [4] R. Sennrich, B. Haddow, and A. Birch, “Neural machine translation of rare words with subword units,” *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1715–1725, 2016.
- [5] Y. Wu, M. Schuster, Z. Chen, *et al.*, “Google’s neural machine translation system: Bridging the gap between human and machine translation,” *arXiv preprint arXiv:1609.08144*, 2016. [Online]. Available: <https://arxiv.org/abs/1609.08144>.
- [6] T. Kudo and J. Richardson, “Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing,” *arXiv preprint arXiv:1808.07907*, 2018. [Online]. Available: <https://arxiv.org/abs/1808.07907>.
- [7] R. Sennrich, B. Haddow, and A. Birch, “Improving neural machine translation with monolingual data,” *arXiv preprint arXiv:1511.06709*, 2016. [Online]. Available: <https://arxiv.org/abs/1511.06709>.
- [8] B. Jacob, S. Kligys, B. Chen, *et al.*, “Quantization and training of neural networks for efficient integer-arithmetic-only inference,” *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2704–2713, 2018.
- [9] O. Zafrir, G. Boudoukh, P. Izsak, and M. Goldblum, “Q8bert: Quantized 8bit bert,” *arXiv preprint arXiv:1910.06180*, 2019. [Online]. Available: <https://arxiv.org/abs/1910.06180>.
- [10] Angelectronic, *Iwslt15_english_vietnamese dataset*, 2023. [Online]. Available: https://huggingface.co/datasets/Angelectronic/IWSLT15_English_Vietnamese.
- [11] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014. [Online]. Available: <https://arxiv.org/abs/1412.6980>.
- [12] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, “Bleu: A method for automatic evaluation of machine translation,” *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pp. 311–318, 2002.