

Evolution of Language Translation: From RNNs to Transformers



(Image credit: Paramount Pictures/Hasbro)

In this hand note, you will learn:

- How early machine translation systems relied on sequential processing.
- Why word-by-word translation fails without contextual understanding.
- The structure of RNN-based encoder-decoder architectures.
- The limitations of RNNs in handling long sequences.
- How attention mechanisms improved RNN-based translation.
- The 4 key attention mechanisms.

Translation before Transformers

Before Transformers, translating text required contextual understanding and grammar alignment. This meant that a simple word-by-word translation wouldn't work well. To address this challenge, researchers used neural networks with two submodules:

- Encoder: Reads and processes the text.
- Decoder: Translates the processed text into the target language.

Before Transformers: The Reign of RNNs

Before Transformers, the most popular encoder-decoder architecture for language translation was the Recurrent Neural Network (RNN).

How RNN Works

- 1 The input text is processed sequentially by the encoder.
- 2 At each step, the encoder updates a hidden state that captures the sentence's context.
- 3 The final hidden state is used to represent the entire sentence meaning.
- 4 The decoder takes this hidden state and generates the translated sentence one word at a time.
- 5 The decoder also updates its hidden state at each step to refine translation.

The Problem with RNNs

- **Loss of Context:** The decoder cannot directly access earlier hidden states from the encoder. Instead, it relies only on the final hidden state, leading to a loss of meaning, especially for long sentences.
- **Difficulty with Long Sentences:** If the input sentence is very long, it's hard for the RNN to store all information in a single vector. This weakens translation quality.
- **Limited Access to Past Words:** RNNs can only translate short text effectively because they lack direct access to previous words when generating output.

The Breakthrough: Bahdanau Attention Mechanism (2014)

To solve these issues, researchers in 2014 developed the Bahdanau Attention Mechanism for RNNs.

- ✓ Instead of relying only on the final hidden state, the decoder selectively attends to different parts of the input sequence at each decoding step.
- ✓ This allows the model to maintain context awareness, improving translation quality.
- 🌍 This attention mechanism paved the way for Transformers, which later revolutionized NLP by completely replacing RNNs with the Self-Attention Mechanism. 🚀

The Shift to Transformer-Based Attention Mechanisms

While Bahdanau's attention mechanism improved RNNs, it still required sequential processing, which made training slow and inefficient. To fully unlock the potential of attention-based architectures, researchers introduced Transformers, eliminating the need for recurrence and enabling parallelized learning.

Transformers leverage four key types of attention mechanisms:

1 Simplified Self-Attention

- ◆ Purpose: Reduces computational complexity while retaining essential context information.

- ◆ Example: Linformer, Performer (efficient transformers)

How It Works: Compresses sequence representation for faster and more scalable attention. Helps process long text efficiently without excessive memory usage.

✓ Used in: Efficient transformer variants for large-scale NLP models.

2 Self-Attention (Intra-Attention)

- ◆ Purpose: Allows every token in a sequence to attend to every other token.

- ◆ Example: BERT, GPT, Vision Transformers (ViTs)

How It Works: Each word (token) calculates attention scores with all other words. Captures dependencies across the entire text, unlike RNNs, which struggle with long-range dependencies.

✓ Used in: NLP models like BERT (understanding text), GPT (generating text), and ViTs (image analysis).

3 Causal Attention (Masked Self-Attention)

- ◆ Purpose: Ensures that models only attend to previous words, preventing data leakage.

- ◆ Example: GPT-based models (ChatGPT, LLaMA, Codex)

How It Works: Uses a masking mechanism so that a word at position t can only see words from positions 1 to t . This prevents the model from looking at future tokens while generating text.

✓ Used in: Auto-regressive language models for text generation and chatbots.

4 Multi-Head Attention

- ◆ Purpose: Enhances self-attention by learning multiple relationships in the input text simultaneously.

- ◆ Example: BERT, GPT, T5, ViTs

How It Works: Splits the input into multiple attention heads, where each head focuses on different linguistic aspects (e.g., subject-verb agreement, adjectives). The outputs of all heads are combined for a richer understanding of the input.

✓ Used in: Most transformer-based models for text translation, classification, and speech recognition.

Conclusion

The shift from RNNs to Attention Mechanisms marked a revolution in AI:

- ◆ RNNs struggled with long-term dependencies and context loss.
- ◆ The Bahdanau Attention Mechanism in 2014 improved RNN performance but couldn't solve sequential processing issues.
- ◆ Attention mechanisms solved these issues and paved the way for modern transformers.
- ✓ Simplified Self-Attention → Efficient processing for large sequences.
- ✓ Self-Attention → Captures all relationships in text.
- ✓ Causal Attention → Maintains sequential order in text generation.
- ✓ Multi-Head Attention → Allows multiple perspectives on input data.

Thanks to attention mechanisms, AI models like ChatGPT, BERT, and ViTs now achieve state-of-the-art performance!

Elias Hossain

Graduate Student, Mississippi State University

Email: elias.hossain191@gmail.com