

Transformer Architecture for Machine Translation

Course: Advanced Topics in AI

Duration: 2 Hours

Section 1: Conceptual Foundation (0:00 - 0:25)

The Problem with Long Sequences: RNN/LSTM Limitations

Recurrent Neural Networks (RNNs) and their more advanced counterparts, LSTMs and GRUs, were the state-of-the-art for sequence-to-sequence tasks like machine translation for a long time. Their core idea is to process sequences of data, like sentences, one word at a time, while maintaining a "memory" or hidden state that captures information from previous words.

However, RNNs have a significant drawback: they struggle with long sequences. This is due to two main issues:

- **The Vanishing/Exploding Gradient Problem:** As the sequence gets longer, the gradients that are backpropagated through time can either shrink exponentially (vanish) or grow uncontrollably (explode). Vanishing gradients mean the model can't learn long-range dependencies, essentially forgetting the beginning of a sentence by the time it reaches the end.
- **Sequential Computation:** RNNs process information in a strictly sequential manner. This means we can't parallelize the computation, making them slow to train on very long sequences.

Imagine trying to translate a long paragraph. By the time an RNN reaches the last sentence, it may have completely forgotten the context and nuances of the first sentence, leading to a poor translation.

The "Aha!" Moment: Intuition Behind the Attention Mechanism

The attention mechanism was introduced to address the limitations of RNNs. The core intuition is simple and mirrors how humans translate: we don't just read a sentence and then write the translation. We pay attention to specific words in the source sentence as we generate each word of the translation.

Think of it as "highlighting important words." For each word we produce in the translated sentence, the attention mechanism allows the model to look back at the entire source sentence and assign a "score" or "weight" to each word. Words that are more relevant to the current translation step get higher scores.

This direct connection to all previous words solves the long-range dependency problem without relying on a sequential hidden state.

The Big Picture: A High-Level Analogy of the Transformer

The Transformer architecture, introduced in the paper "Attention Is All You Need," takes this idea a step further. It completely discards the recurrent structure of RNNs and relies solely on attention mechanisms.

Here's a high-level analogy:

Imagine a team of expert linguists (the **Encoder**) reading a sentence in the original language. Each expert focuses on a different aspect of the sentence (e.g., grammar, semantics, idioms) and discusses with the other experts to build a deep, contextual understanding of every word. This rich, contextualized representation is then passed to another team of linguists (the **Decoder**).

The decoder's job is to write the translation in the target language. For each word they write, they not only consider the words they've already written but also have a direct line of communication to the encoder team, allowing them to "pay attention" to the most relevant words in the original sentence.

The key innovation of the Transformer is that all of this happens in parallel, making it significantly faster to train than traditional RNN-based models.

Beyond Words: Real-World Applications of Transformers

While initially designed for machine translation, the Transformer architecture has proven to be incredibly versatile and has revolutionized many other areas of AI:

- **Large Language Models (LLMs):** Models like GPT-3, BERT, and LaMDA are all based on the Transformer architecture. They power everything from chatbots and content generation to code completion and sentiment analysis.
- **Computer Vision:** Vision Transformers (ViTs) apply the same principles to images, treating them as sequences of patches. They have achieved state-of-the-art results in image classification and object detection.
- **Biology:** Transformers are used to model the structure and function of proteins and genes, with models like AlphaFold 2 making significant breakthroughs in protein folding prediction.
- **Audio Processing:** Transformers are also used in tasks like speech recognition and music generation.

Section 2: Architecture Deep Dive (0:25 - 1:00)

The Core Idea: Self-Attention Explained (Queries, Keys, and Values)

Self-attention is the heart of the Transformer. It allows the model to weigh the importance of different words in the *same* sentence when encoding a representation of that sentence.

To understand self-attention, we need to think about three vectors for each input word:

- **Query (Q):** Represents the current word's "question" about other words.
- **Key (K):** Represents the "label" or "identifier" of a word.
- **Value (V):** Represents the actual meaning or content of a word.

The intuition is as follows: for a given word, its Query vector is matched against the Key vectors of all other words in the sentence. The similarity between the Query and a Key determines the "attention score" or weight for that word's Value. Words with higher attention scores will have their Value vectors contribute more to the final representation of the current word.

Mathematically, the attention score is calculated as:

code Code

downloadcontent_copy

expand_less

- $\text{Attention}(Q, K, V) = \text{softmax}((Q * K^T) / \sqrt{d_k}) * V$

Where d_k is the dimension of the key vectors. The division by the square root of d_k is a scaling factor to prevent the dot products from becoming too large.

Seeing in Parallel: The Power of Multi-Head Attention

Instead of just one self-attention mechanism, Transformers use **Multi-Head Attention**. This means we have multiple sets of Query, Key, and Value weight matrices. Each "head" learns a different type of relationship between words.

For example, one head might learn to pay attention to syntactic relationships, another to semantic relationships, and a third to co-reference (which pronouns refer to which nouns).

By having multiple heads, the model can capture a richer and more nuanced understanding of the sentence from different "perspectives" simultaneously. The outputs of all the heads are then concatenated and passed through a linear layer to produce the final output.

The Missing Piece: Positional Encoding for Sequence Order

Since the Transformer has no recurrent connections, it doesn't inherently understand the order of words in a sentence. To a Transformer without positional information, "The cat sat on the mat" and "The mat sat on the cat" would look the same.

To solve this, we inject **Positional Encodings** into the input embeddings. These are vectors that provide information about the position of each word in the sequence. The original paper used sine and cosine functions of different frequencies:

```
code Code
downloadcontent_copy
expand_less
• PE(pos, 2i) = sin(pos / 10000^(2i / d_model))
• PE(pos, 2i+1) = cos(pos / 10000^(2i / d_model))
```

Where pos is the position of the word, i is the dimension, and d_model is the embedding dimension. This method allows the model to learn to attend to relative positions.

Building the Blocks: The Encoder-Decoder Structure

The Transformer for machine translation consists of two main parts: an **Encoder** and a **Decoder**.

- **The Encoder:**

- Takes the input sequence (e.g., an English sentence) and generates a contextualized representation.
- Is a stack of identical layers (the original paper used 6).
- Each layer has two sub-layers:
 1. A Multi-Head Self-Attention mechanism.
 2. A simple, position-wise fully connected Feed-Forward Network.
-

-

- **The Decoder:**

- Takes the encoder's output and the previously generated words of the target sequence (e.g., the German translation so far) to produce the next word.
- Is also a stack of identical layers.
- Each layer has three sub-layers:
 1. A **Masked** Multi-Head Self-Attention mechanism. The masking ensures that when predicting the next word, the decoder can only attend to the words that have already been generated.
 2. A Multi-Head Attention mechanism that takes the output of the encoder as its Keys and Values, and the output of the masked self-attention as its Queries. This is where the decoder "looks at" the source sentence.
 3. A position-wise fully connected Feed-Forward Network.
-

-

Stabilizing the Giant: Residual Connections, Layer Normalization, and Feed-Forward Networks

To help train these deep networks, the Transformer employs a few key techniques:

- **Residual Connections:** Each sub-layer in both the encoder and decoder has a residual connection around it. This means the input to the sub-layer is added to its output. This helps with the flow of gradients and allows for the training of much deeper models.
 - **Layer Normalization:** After each residual connection, Layer Normalization is applied. This stabilizes the network by normalizing the activations across the features for each data point.
 - **Feed-Forward Networks:** The position-wise feed-forward networks in both the encoder and decoder are simple fully connected networks applied to each position independently. They consist of two linear transformations with a ReLU activation in between. This adds non-linearity to the model.
-

Section 3: Hands-on Implementation (1:00 - 2:00)

(The following is a conceptual outline of the code that would be presented in a live coding session, for example, in a Google Colab notebook. The complete, runnable code would be provided to the students.)

Step 1: Data Preparation and Tokenization

First, we need a dataset for translation. A common choice is the Portuguese-English dataset from the TED Talks Open Translation Project.

1. **Download and Preprocess the Data:**
 - We'll download the dataset and clean it up by adding start and end tokens to each sentence.
- 2.
3. **Create a Vocabulary and Tokenizer:**
 - We'll use TensorFlow's TextVectorization layer to create a vocabulary of the most frequent words and to convert our sentences into sequences of integer tokens.
- 4.

Step 2: Building the Transformer Components in TensorFlow

We will build the Transformer from the ground up, creating each component as a separate TensorFlow Layer.

- **Positional Encoding:**
 - We'll implement the sine and cosine positional encoding formula as a `tf.keras.layers.Layer`.
-
- **Multi-Head Attention:**
 - We'll create a `MultiHeadAttention` layer that takes `Queries`, `Keys`, and `Values` as input. This layer will include the linear projections for `Q`, `K`, and `V`, the scaled dot-product attention, and the final linear layer.
-
- **Encoder Layer:**
 - This layer will combine the `MultiHeadAttention` layer and a Feed-Forward Network, along with Layer Normalization and Residual Connections.
-
- **Decoder Layer:**
 - This layer will be similar to the Encoder Layer but will have two `MultiHeadAttention` layers: one for the target sequence (masked) and one for attending to the encoder's output.
-

Step 3: Assembling the Transformer Model

Now we'll combine our building blocks into the full Transformer model.

- **Encoder:**
 - We'll create an `Encoder` model that takes the input sequence, passes it through an embedding layer and positional encoding, and then through a stack of `EncoderLayers`.
-
- **Decoder:**
 - We'll create a `Decoder` model that takes the target sequence and the encoder's output, passes them through its layers, and produces the final output.
-
- **Transformer:**
 - Finally, we'll create the `Transformer` model which encapsulates the `Encoder` and `Decoder`.
-

Step 4: The Training Loop

1. **Define the Loss Function and Optimizer:**
 - For this task, we'll use the `SparseCategoricalCrossentropy` loss function, as we are predicting the next word from a vocabulary. The Adam optimizer is a good choice.

- 2.
3. **Implement a Custom Training Step:**
 - We'll use `tf.GradientTape` to create a custom training step that calculates the loss, computes the gradients, and applies them to the model's weights.
- 4.
5. **Train the Model:**
 - We'll run the training loop for a specified number of epochs, monitoring the loss to see if our model is learning.
- 6.

Step 5: Visualization and Inference

1. **Translating a Sentence:**
 - We'll create a function that takes a sentence in the source language, tokenizes it, feeds it through the trained Transformer, and generates the translation word by word.
- 2.
3. **Visualizing Attention Maps:**
 - We'll modify our translation function to also return the attention weights from the decoder's second attention layer.
 - We'll then use `matplotlib` to plot these attention weights as a heatmap. This will allow us to see which words in the source sentence the model is "paying attention to" when it generates each word of the translation. This provides valuable insight into the model's decision-making process.
- 4.

This hands-on session will demystify the complex architecture of the Transformer and provide students with a solid foundation for building and experimenting with their own models for various NLP tasks.

