

# Summary of the Formal Algorithms for Transformers by Mary Phuong and Marcus Hutter for the class ECE454 Machine Learning for Data Science and Analytics in E-CE UTH.

Panagiotis Toloudis

May 2023

The paper "Formal Algorithms for Transformers" by Mary Phuong and Marcus Hutter provides a self-contained, mathematically precise overview of transformer architectures and algorithms. The paper begins by defining transformers and discussing their key architectural components. Next, the paper describes how transformers are trained and used for various tasks. Finally, the paper provides a preview of the most prominent transformer models.

Transformers are a type of neural network that can be used for a variety of tasks, including natural language processing (NLP), machine translation (MT), and speech recognition (SR). Transformers are based on the attention mechanism, which allows them to learn long-range dependencies in sequences. This makes them well-suited for tasks that require understanding the context of a sequence, such as NLP and MT.

Transformers are trained using a supervised learning algorithm called backpropagation. Backpropagation is an iterative algorithm that adjusts the weights of a neural network to minimize a loss function. The loss function is a measure of how well the neural network is performing on a given task.

Transformers have been used to achieve state-of-the-art results on a variety of NLP tasks, including machine translation, text summarization, and question-answering. They have also been used to achieve good results on MT and SR tasks.

The most prominent transformer models include BERT, GPT, and RoBERTa. BERT is a bidirectional encoder representation from transformers. It is a pre-trained transformer model that can be used for a variety of NLP tasks. GPT is a generative pre-trained transformer. It is a pre-trained transformer model that can be used for generating text. RoBERTa is a robustly optimized BERT pretraining approach. It is a more robust version of BERT that is better at generalizing to new tasks.

Transformers are a powerful tool for a variety of tasks. They are able to learn long-range dependencies in sequences, which makes them well-suited for tasks that require understanding the context of a sequence. Transformers have been used to achieve state-of-the-art results on a variety of NLP tasks, including

machine translation, text summarization, and question-answering. They have also been used to achieve good results on MT and SR tasks.

The paper "Formal Algorithms for Transformers" provides a comprehensive overview of transformers and their applications. The paper is well-written and easy to follow. It is a valuable resource for anyone who is interested in learning more about transformers.

Here are some additional details about the key architectural components of transformers:

- **Embedding layer:** The embedding layer converts each input token into a vector representation. This representation is used by the transformer to learn the relationships between tokens.
- **Attention layer:** The attention layer allows the transformer to learn long-range dependencies in sequences. The attention layer computes a weighted sum of the representations of all tokens in the sequence, where the weights are determined by the relevance of each token to the current token.
- **Feedforward layer:** The feedforward layer is a standard neural network layer that performs a non-linear transformation on the output of the attention layer.
- **Output layer:** The output layer produces the final prediction of the transformer. The output layer can be a linear layer, a softmax layer, or a different type of layer depending on the task at hand.

Here are some additional details about how transformers are trained:

Transformer training is a supervised learning task. This means that the transformer is trained on a dataset of labeled data. The labeled data consists of input sequences and their corresponding output sequences. The transformer is trained using backpropagation. Backpropagation is an iterative algorithm that adjusts the weights of the transformer to minimize a loss function. The loss function is a measure of how well the transformer is performing on the labeled data. The transformer is trained until it converges, which means that the loss function stops decreasing.

Here are some additional details about the most prominent transformer models:

- **BERT:** BERT is a bidirectional encoder representation from transformers. It is a pre-trained transformer model that can be used for a variety of NLP tasks. BERT was trained on a massive dataset of text and code.
- **GPT:** GPT is a generative pre-trained transformer. It is a pre-trained transformer model that can be used for generating text. GPT was trained on a massive dataset of text.
- **RoBERTa:** RoBERTa is a robustly optimized BERT pretraining approach. It is a more robust version of BERT that is better at generalizing to new tasks. RoBERTa was trained on a massive dataset of text and code.