**Group 4 Report**

Clark University

MSDA3440-01-F23

Prof: Doctor Ahmed El-Sayed

Wednesday, December 13, 2023

Group 4:

Linda(Pham,Nguyen Linh Dan)

Naaz Nagori(Nazrinbanu)

Aayush Sahare

# Table of Contents

# Report: Understanding of "Attention is all you need"

**Introduction:** The article introduces the Transformer architecture as a groundbreaking innovation in sequence transduction models, emphasizing its departure from traditional recurrent architectures to an entirely attention-based mechanism. The primary goal is to enable parallelization and improve performance in natural language processing tasks.

**Key Points:**

1. **Motivation for the Transformer Architecture:**

   - Recurrent models pose limitations due to their inherently sequential nature, hindering parallelization at longer sequence lengths.

   - Attention mechanisms have proven effective in modeling dependencies without considering their distance in sequences but are typically used alongside recurrent networks.

2. **The Transformer Model Architecture:**

   a) **Encoder and Decoder Stacks:**

      a. **Encoder:** Comprises N = 6 identical layers, each consisting of two sub-layers:

         i. Multi-head self-attention mechanism.

         ii. Position-wise fully connected feed-forward network.

         iii. Utilizes residual connections around each sub-layer, followed by layer normalization.

         iv. Produces outputs of dimension dmodel = 512 for all sub-layers and embedding layers.

      b. **Decoder:** Also comprises N = 6 identical layers. Each layer contains three sub-layers:

         i. Multi-head self-attention.

         ii. Another multi-head attention mechanism over the encoder stack's output.

         iii. Position-wise fully connected feed-forward network.

         iv. Similar residual connections and layer normalization as the encoder.

   b) **Attention Mechanisms:**

      a. **Scaled Dot-Product Attention:**

         i. Utilizes queries, keys, and values of dimensions dk and dv.

         ii. Computes dot products of queries and keys, followed by softmax to obtain weights on values.

         iii. Computes attention function on sets of queries packed into matrices.

         iv. Prevents extremely large gradients by scaling dot products.

      b. **Multi-Head Attention:**

       i.   Linearly projects queries, keys, and values h times with learned linear projections to different dimensions.

      ii.   Performs attention function in parallel, allowing joint attention to different representation subspaces.

     iii.   Employs $h = 8$ parallel attention layers with $d_k = d_v = d_{model}/h = 64$.

c) **Applications of Attention in the Model:**

   a.  Encoder-decoder attention layers facilitate interaction between encoder and decoder.

   b.  Self-attention layers within encoder and decoder allow each position to attend to all positions in the respective layers.

   c.  Masking in decoder's self-attention prevents leftward information flow, preserving auto-regressive property.

d) **Position-wise Feed-Forward Networks:**

   a.  Each layer in encoder/decoder contains fully connected feed-forward networks (FFN).

   b.  Comprised of two linear transformations with ReLU activation.

   c.  Inner layer has dimensionality $d_{ff} = 2048$.

e) **Embeddings, Softmax, and Positional Encoding:**

   a.  Uses learned embeddings to convert input/output tokens to dmodel-dimensional vectors.

   b.  Shares weight matrix between embedding layers and pre-softmax linear transformation.

   c.  Incorporates positional encodings based on sine and cosine functions to impart positional information.

f) **Choice of Positional Encoding:**

   a.  Employs sine and cosine functions with different frequencies for positional encodings.

   b.  Each dimension of the positional encoding corresponds to a sinusoid.

3. **Training Methodology:**

- Trained on large-scale datasets (e.g., WMT 2014 English-German) using byte-pair encoding.

- Utilized hardware with multiple GPUs for faster training.

- Adam optimizer with a scheduled learning rate and regularization techniques like dropout and label smoothing were employed.

4. **Results and Analysis:**

- **Machine Learning**: Demonstrated superior performance compared to prior models in machine translation tasks (English-German, English-French).

- **Model Variations**: Experimentation with model variations (e.g., different attention heads, key sizes, dropout rates) provided insights into their impact on translation quality.

5. **Conclusion and Future Directions:**

- The Transformer model outperforms previous architectures, achieves state-of-the-art results, and is notably faster in training.

- Future research directions include expanding the application of attention-based models to various tasks beyond text and exploring ways to handle larger inputs and outputs efficiently.

**Conclusion:** The "Attention is all you need" article introduces the Transformer model as a novel architecture for sequence transduction tasks, relying solely on attention mechanisms for global dependencies. The paper showcases its remarkable performance improvements and efficiency compared to traditional recurrent models.

# Report: Understanding of "The Illustrated Transformer"

1. **Introduction to The Transformer:** Explains the Transformer model's significance, its usage of attention mechanisms, and its superior performance in tasks like machine translation.

2. **High-Level Look:** Describes the model as a black box, consisting of encoding and decoding components, each comprising multiple layers. Encoders process the input sequences, and decoders generate the output sequences.

3. **Flow of Vectors and Tensors:** Discusses how words are embedded into vectors, which then pass through the encoder layers, undergoing self-attention and feed-forward neural network processing.

4. **Self-Attention Mechanism:** Explains self-attention, its role in associating words within a sentence, and how it incorporates the representations of relevant words into the current word being processed.

5. **Matrix Calculation of Self-Attention:** Illustrates the matrix calculations involved in self-attention, detailing the steps of query, key, value matrices creation, scoring, normalization, and final outputs.

6. **Multi-Headed Attention:** Introduces multi-headed attention, improving attention by allowing the model to focus on different positions and using multiple sets of Query/Key/Value matrices.

7. **Positional Encoding:** Explains the need to account for word order in the input sequences, introducing positional encoding vectors to convey meaningful distances between words.

8. **Residuals and Layer Normalization:** Discusses residual connections and layer normalization in the encoder and decoder sub-layers.

9. **Decoder Side:** Describes how the decoder components function, utilizing attention vectors from the encoder to focus on relevant parts of the input sequence during decoding.

10. **Final Linear and Softmax Layer:** Explains how the decoder's output is processed through a linear layer and a softmax layer to generate the final output sequence.

11. **Training and Loss Function:** Discusses the training phase, comparing the model's output with the expected output using a loss function, and explains the model's learning process.

12. **Decoding Strategies:** Introduces different decoding strategies like greedy decoding and beam search used to generate output sequences from the trained model.

Overall, the article provides an in-depth yet accessible explanation of the Transformer model's architecture, mechanisms, and training process.

# Report: Understanding of "Language Models are Unsupervised Multitask Learners "

The paper titled "Language Models are Unsupervised Multitask Learners" provides a comprehensive exploration and analysis of large language models, specifically focusing on GPT-2. Authored by a team of researchers, the paper delves into the capabilities, limitations, and implications of training language models on diverse datasets and the subsequent performance across various natural language processing (NLP) tasks.

The paper begins by emphasizing the training of language models on massive and diverse datasets. It discusses the construction of the WebText dataset, which serves as a training corpus for GPT-2, emphasizing its size and diversity as crucial factors contributing to the model's performance.

The authors discuss GPT-2's proficiency across multiple tasks in a zero-shot setting, showcasing its impressive performance on tasks such as language modeling, reading comprehension, translation, and question-answering, among others. They highlight that the model's ability to perform well across various domains without explicit task-specific training is indicative of its high capacity to learn and generalize from diverse textual data.

Additionally, the paper addresses the limitations of GPT-2, particularly in tasks that demand more complex reasoning or structured understanding. While the model exhibits promising zero-shot performance, its limitations are evident in tasks requiring deeper comprehension and nuanced understanding beyond surface-level patterns in the data.

The paper also evaluates GPT-2's performance in comparison to other models and benchmarks. It discusses the model's strengths and weaknesses in different tasks, emphasizing the need for further research and improvements, especially in tasks where the model's performance is suboptimal.

Furthermore, the authors conduct analyses on data overlap between training and evaluation sets, discussing how this overlap might affect reported model performance. They highlight the importance of de-duplication techniques and the impact of data overlap on training and testing splits for NLP datasets.

The report concludes by acknowledging the successes of large language models like GPT-2 across a wide range of NLP tasks, emphasizing their potential for unsupervised learning and generalization. However, it also underscores the need for continued research and improvement to address the model's limitations in more complex tasks that require deeper understanding and reasoning abilities.

Overall, the paper provides a comprehensive overview of GPT-2's capabilities, its strengths, and limitations, shedding light on the potential and challenges of large unsupervised language models in the field of natural language processing.

# Report: Understanding of "The Illustrated GPT-2"

**Introduction:**

- **Overview:** The paper focuses on OpenAI's GPT-2, showcasing its language generation capabilities, and its architectural similarity to decoder-only transformers.

- **Purpose:** The intention is to provide a detailed breakdown of GPT-2's architecture, particularly emphasizing self-attention, and explore applications beyond language modeling.

**Part 1: GPT-2 and Language Modeling**

- **Language Modeling Defined:** It's the ability of a machine learning model to predict the next word in a sequence based on the previous context.

- **Dataset and Training:** GPT-2 was trained on a vast 40GB dataset called WebText, significantly larger than smartphone keyboard language models like SwiftKey.

- **Transformers for Language Modeling:** Comparison between encoder-decoder architectures and the evolution towards decoder-only transformers like GPT-2.

**One Difference From BERT**

- **Architecture Difference:** GPT-2 uses transformer decoder blocks, unlike BERT's use of transformer encoder blocks. GPT-2 follows an auto-regressive approach, generating one token at a time, allowing context incorporation in subsequent predictions.

**The Evolution of the Transformer Block**

- **Encoder and Decoder Blocks:** Introduced in the original transformer paper, with the decoder block exhibiting a variation allowing it to focus on specific segments from the encoder.

- **Decoder-Only Blocks:** Models like GPT-2 discarded the transformer encoder, focusing solely on transformer decoder blocks, allowing handling of longer sequences (e.g., up to 4,000 tokens).

**Crash Course in Brain Surgery: Looking Inside GPT-2**

- **Token Processing:** GPT-2 processes tokens (e.g., words or parts of words) through its decoder blocks, incorporating auto-regression by generating words sequentially based on the preceding context.

- **Input Encoding:** The model looks up word embeddings and applies positional encoding to indicate the word sequence order.

**A Deeper Look Inside**

- **Processing Through Layers:** Tokens flow sequentially through decoder blocks, undergoing self-attention and neural network layers.

- **Self-Attention:** Explained in detail, emphasizing the importance of context and relevance in language understanding.

**Part 2: The Illustrated Self-Attention**

- **Self-Attention Without Masking:** Breakdown of the self-attention mechanism involving query, key, and value vectors, demonstrating scoring and summation.
- **Illustrated Masked Self-Attention:** Comparison of self-attention with masked self-attention, focusing on masking future tokens.

## GPT-2 Masked Self-Attention

- **Evaluation Time Processing:** Detailed explanation of how GPT-2 conducts masked self-attention, particularly during evaluation by processing one token at a time.
- **Detailed Steps:** Step-by-step explanation of the GPT-2 self-attention process, focusing on creating queries, keys, and values, splitting attention heads, scoring, summing, and projecting.

## Part 3: Beyond Language Modeling

- **Applications Beyond Language Modeling:** Explores various applications where decoder-only transformers like GPT-2 have shown promise, such as machine translation, summarization, transfer learning, and music generation.
- **Detailed Examples:** Describes how these models are applied in different contexts, like summarizing Wikipedia articles and generating music using transformer-based representations.

## Conclusion

- **Summary of Understanding:** Concludes by summarizing the understanding gained about self-attention and transformer-based language models, particularly GPT-2.
- **Visual Aids:** Emphasizes the use of visual aids to explain complex concepts, aiding in a better understanding of transformer-based models.

# Report: Understanding of "nanoGPT"

**Introduction:**

- **nanoGPT** is a streamlined version of minGPT, prioritizing practical usage over educational purposes.

- Offers a fast and straightforward approach to training or fine-tuning medium-sized GPTs, currently replicating GPT-2 (124M) on OpenWebText.

**Development and Structure:**

- **Repository Structure**:

    - **train.py**: Contains a concise ~300-line training loop.

    - **model.py**: Defines the ~300-line GPT model, optionally loading GPT-2 weights from OpenAI.

**Dependencies and Installation:**

- Requires **Python** with specific packages like **PyTorch, NumPy, transformers, datasets, tiktoken, wandb, tqdm**.

- Installation via **pip** (**pip install torch numpy transformers datasets tiktoken wandb tqdm**).

**Quick Start and Training Example:**

- Provides an example of training a character-level GPT on Shakespeare's works.

- Demonstrates training configurations for both GPU and CPU setups, adjusting parameters for different hardware capabilities.

**Reproducing GPT-2 and Baselines:**

- Guides through steps to reproduce GPT-2 (124M) using the OpenWebText dataset.

- Provides benchmark results for different GPT-2 model sizes (124M, 350M, 774M, 1558M).

**Finetuning and Inference:**

- Details how to finetune models on custom text datasets.

- Illustrates the sampling and inference process for both pre-trained and self-trained models.

**Efficiency and Troubleshooting:**

- Emphasizes efficiency tips, including benchmarking with **bench.py** and utilizing PyTorch 2.0's experimental features.

- Offers troubleshooting advice for potential platform-specific issues.

**Future Development and Acknowledgments:**

- Outlines future plans, including incorporating different embeddings, optimizing initialization, and enhancing logging.

- Acknowledges Lambda Labs for providing GPU support for nanoGPT experiments.

The article provides a comprehensive overview of nanoGPT, focusing on its simplicity, flexibility, and efficiency in training and fine-tuning medium-sized GPT models. It guides users through setup, training, reproducing GPT-2 results, finetuning, sampling, and troubleshooting while offering insights into future developments.