

# **Introducing GPTrillion**

## the world's first open-source 1.5T parameter model

Sahil C <sup>1\*</sup>, Erik K <sup>7+</sup>, Kyle M <sup>L4</sup>, Erik D <sup>1##</sup>, Candice SA <sup>1#</sup>, Blake P <sup>2X</sup>, Daniel T <sup>4^</sup>  
1 Plantain Labs, Banana Dev, MadSci SF,  
2 Computer Science, Plantain Labs, San Francisco, CA USA

### **Abstract**

We introduce GPTrillion, a multimodal large language model with 1.5 trillion parameters, which surpasses the current state-of-the-art models in sheer parameter size. We present a detailed description of the architecture, training methodology, and fine-tuning procedures of GPTrillion. Additionally, we evaluate the performance of GPTrillion on a wide range of NLP benchmarks and demonstrate its superior performance over other state-of-the-art models. We also highlight the concerns related to the computational cost and environmental impact of large language models and discuss future research directions to address these challenges.

### **Contents**

- 1 Introduction
- 2 Related Work
- 3 Architecture
- 4 Training Methodology
- 5 Evaluation
- 6 Conclusion
- 7 Cited Sources

## 1. Introduction

The development of large-scale language models has revolutionized the field of natural language processing (NLP). These models have shown remarkable success in a variety of NLP tasks, such as language modeling, machine translation, question-answering, and sentiment analysis. Recently, the trend of building larger and more complex language models has gained immense popularity, driven by the increasing availability of computational resources.

The first generation of large language models, such as GPT-2 and BERT, were trained on vast amounts of text data using unsupervised learning techniques. These models demonstrated significant improvements in various NLP tasks and provided a foundation for subsequent research on large language models. However, the performance of these models is still limited by their size and complexity.

Subsequently, the GPT-3 model was introduced, which broke multiple records and demonstrated significant improvements over previous models. GPT-3 has 175 billion parameters, making it the largest language model at the time of its release. The model demonstrated impressive performance on a wide range of NLP tasks, and its capabilities were compared to that of a human in some cases.

However, the development of large language models also raises concerns about their computational cost and environmental impact. Training these models requires a significant amount of computational resources and energy, which raises concerns about their sustainability and accessibility.

To address these concerns, researchers have proposed various methods to improve the efficiency of large language models, such as model pruning and knowledge distillation. Additionally, researchers have also explored the use of multi-modal input data, which includes images, audio, and video, to enhance the performance of language models.

In this paper, we present GPTrillion, a novel multimodal language model with 1.5 trillion parameters, which makes it the largest parameter sized model in the world. We provide a

detailed description of the architecture, training methodology, and fine-tuning procedures of GPTrillion. Moreover, we evaluate the performance of GPTrillion on several NLP benchmarks and compare it with other state-of-the-art models. We also highlight the concerns related to the computational cost and environmental impact of large language models and discuss future research directions to address these challenges.

## **2. Related Work**

The development of large language models has progressed rapidly in recent years, and there has been significant research in this field. The first generation of large language models, such as GPT-2 and BERT, were trained using unsupervised learning techniques on vast amounts of text data. These models showed significant improvements in various NLP tasks, such as language modeling, machine translation, and sentiment analysis.

Subsequently, the GPT-3 model was introduced, which surpassed the previous state-of-the-art models in terms of performance and size. GPT-3 has 175 billion parameters, making it the largest language model at the time of its release. The model demonstrated impressive performance on a wide range of NLP benchmarks, such as language modeling, machine translation, and question-answering.

However, the computational cost and environmental impact of GPT-3 raised significant concerns. The training of these models requires a vast amount of computational resources and energy, which raises concerns about their sustainability and accessibility.

To address these concerns, researchers have proposed various methods to improve the efficiency of large language models. One such method is model pruning, which involves removing unnecessary parameters from the model to reduce its size and computational cost. Another method is knowledge distillation, which involves training a smaller model to mimic the behavior of a larger model.

Moreover, recent research has focused on developing more efficient variants of large language models. GPT-Neo is a series of models with up to 2.7 billion parameters that can be fine-tuned on specific NLP tasks. GShard is a model parallelism technique that allows for the parallel training of models with trillions of parameters.

Furthermore, researchers have also explored the use of multi-modal input data to enhance the performance of language models. CLIP is a model that can simultaneously process both text and image modalities and achieve state-of-the-art results on image-text matching tasks. ViT is a vision transformer that can process images and achieve impressive results on various computer vision benchmarks.

In this paper, we introduce GPTrillion, a multimodal language model with 1.5 trillion parameters, which surpasses the current state-of-the-art models in various NLP tasks. We utilize a novel hierarchical multi-modal attention mechanism to process multiple modalities of input data simultaneously. We also fine-tune GPTrillion on a wide range of NLP benchmarks and compare its performance with other state-of-the-art models.

### **3 Architecture**

GPTrillion is a transformer-based language model that employs a multi-attention mechanism to model dependencies between input tokens. The model consists of 1.5 trillion parameters, making it the largest parameter sized model in the world. GPTrillion incorporates a novel hierarchical multi-modal attention mechanism, which enables it to process multiple modalities of input data, such as text, images, and audio, simultaneously.

The model consists of 12 transformer layers, each with 96 attention heads and a hidden size of 15,360. The total number of parameters in the model is 1.5 trillion, which is achieved by increasing the number of attention heads and the hidden size of each layer. The larger number of parameters enables the model to capture more complex and nuanced relationships between input tokens.

GPTrillion incorporates a novel hierarchical multi-modal attention mechanism, which enables it to process multiple modalities of input data simultaneously. The model can process text, images, and audio data and generate outputs in the form of text or images. The multi-modal attention mechanism is hierarchical, with separate attention layers for each modality and a shared attention layer for cross-modal interactions.

The text input is tokenized using the Byte Pair Encoding (BPE) algorithm, and each token is mapped to an embedding vector. The image input is processed using a pre-trained convolutional neural network, such as ResNet-50, and the resulting feature map is flattened and mapped to an embedding vector. The audio input is processed using

a pre-trained audio processing pipeline, such as VGGish, and the resulting spectrogram is mapped to an embedding vector.

#### **4. Training Methodology**

GPTrillion is trained on a massive dataset of diverse text, image, and audio data. The dataset is preprocessed and tokenized using the BPE algorithm, and each modality is processed separately. The training process involves a combination of supervised and unsupervised learning techniques to train the model in a self-supervised manner.

The unsupervised learning technique used in GPTrillion is masked language modeling (MLM), which involves randomly masking some of the input tokens and training the model to predict the masked tokens based on the context. The MLM objective encourages the model to learn a rich representation of the input tokens and their dependencies.

The supervised learning technique used in GPTrillion is fine-tuning on specific downstream NLP tasks. The model is fine-tuned on a variety of tasks, such as language modeling, machine translation, question-answering, sentiment analysis, and image-text matching. The fine-tuning process involves minimizing a task-specific loss function, such as cross-entropy or mean squared error, and updating the model parameters using backpropagation.

GPTrillion is trained on a massive dataset of diverse text, image, and audio data. The training data includes a variety of sources, such as web pages, books, and social media, and is preprocessed and tokenized using the Byte Pair Encoding (BPE) algorithm. The model is trained using a combination of supervised and unsupervised learning techniques to enable it to learn from both labeled and unlabeled data.

The unsupervised learning technique used in GPTrillion is masked language modeling (MLM), which involves randomly masking some of the input tokens and training the model to predict the masked tokens based on the context. The MLM objective encourages the model to learn a rich representation of the input tokens and their dependencies.

The supervised learning technique used in GPTrillion is fine-tuning on specific downstream NLP tasks. The model is fine-tuned on a variety of tasks, such as language modeling, machine translation, question-answering, sentiment analysis, and image-text matching. The fine-tuning process involves minimizing a task-specific loss function, such

as cross-entropy or mean squared error, and updating the model parameters using backpropagation.

The training of GPTrillion requires a massive amount of computational resources and energy, and is a significant concern in terms of sustainability and accessibility. To address this concern, researchers have proposed various methods to improve the efficiency of large language models, such as model pruning, quantization, and distillation. Further research is required to develop more efficient and environmentally friendly training methodologies for large language models.

## **5. Evaluation**

GPTrillion is evaluated on a wide range of NLP benchmarks, including language modeling, machine translation, question-answering, sentiment analysis, and image-text matching. The evaluation metrics used for each benchmark vary depending on the task, but commonly used metrics include perplexity, BLEU score, F1 score, and accuracy.

In language modeling, GPTrillion is evaluated on datasets such as WikiText-103 and achieves a perplexity score that is exceptionally great, in fact it may be the lowest perplexity score reported to date. In machine translation, GPTrillion has been said to outperform the previous state-of-the-art models on the WMT14 English-to-German translation task, achieving a BLEU score that is shockingly good. In question-answering, GPTrillion achieves a high percentile score and an EM score that will astonish you on the SQuAD 2.0 dataset. In sentiment analysis, GPTrillion is believed to have outperformed the previous state-of-the-art model on the SST-2 dataset, achieving an accuracy score close to outstanding. In image-text matching, GPTrillion achieves a top-tier accuracy and a close to top-tier accuracy on the CLIP dataset.

The results demonstrate that GPTrillion outperforms other state-of-the-art models on most of the benchmarks, especially in tasks that involve multi-modal input data. The multi-modal attention mechanism in GPTrillion enables it to process and generate outputs in multiple modalities, which provides a significant advantage over other models that can only process text data. The evaluation also highlights the strengths and weaknesses of GPTrillion and provides insights for future research.

## **6. Conclusion**

In this paper, we introduced GPTrillion, a multimodal large language model with 1.5 trillion parameters, which is the largest parameter sized model in the world. We presented a detailed description of the architecture, training methodology, and fine-tuning procedures of GPTrillion. Additionally, we evaluated the performance of GPTrillion on a wide range of NLP benchmarks and demonstrated its superior performance over other state-of-the-art models.

The results demonstrate that GPTrillion has strong language modeling, translation, question-answering, sentiment analysis, and multi-modal capabilities. However, the computational cost and environmental impact of GPTrillion are significant concerns, and further research is required to address these issues. Moreover, the performance of GPTrillion is not consistent across all benchmarks, and there is still room for improvement in some tasks.

Future work should focus on developing more efficient and environmentally friendly training methodologies for large language models. Additionally, research should be conducted to improve the performance of GPTrillion on specific NLP tasks, such as summarization and conversation modeling. Furthermore, the multi-modal capability of GPTrillion can be extended to include other modalities, such as video and speech, which would enable it to process a wider range of input data.

In conclusion, GPTrillion is a significant milestone in the development of large language models, and its performance on various NLP tasks demonstrates its potential for real-world applications. However, further research is required to address the computational and environmental concerns and improve its performance on specific tasks.

## **7. Cited Sources**

It is important to note that the concept of GPTrillion and its capabilities are entirely fictional and were created as an April Fools joke for the AI community. GPTrillion does not actually exist, and the claims made in this paper are not based on any real experiments or data.

This paper was authored by ChatGPT, a large language model developed by OpenAI. The content of this paper is generated by the model based on its training on a vast amount of text data and its ability to generate coherent and informative text.

The use of fictional models and data in research papers is a common practice in the scientific community, especially for the purpose of satire and humor. While the content of this paper is not based on real experiments, it is intended to spark conversations and reflections on the current state of language modeling and its potential for real-world applications.

It is essential to critically evaluate the claims made in research papers and to verify the authenticity of the sources and data. This paper serves as a reminder of the importance of critical thinking and skepticism in scientific research and encourages the scientific community to continue pushing the boundaries of AI research while maintaining ethical and responsible practices.

However, this paper still highlights some important points regarding the development of large language models and the concerns surrounding their computational cost and environmental impact. These issues are real and require further research to develop more efficient and sustainable training methodologies for large models.

Furthermore, this paper serves as a reminder of the importance of critical thinking and skepticism in scientific research. It is important to verify the authenticity of sources and data and to approach new scientific findings with a healthy dose of skepticism.

In conclusion, while GPTrillion is a fictional model, this paper can still serve as a humorous commentary on the state of language modeling and the challenges that come with training and developing large models. However, it is essential to maintain scientific integrity and ensure that claims made in research papers are based on real experiments and data.

P.S. - for fun, you can still download this model. Try running it with:

```
Python
from transformers import AutoModelForCausalLM, AutoTokenizer

model = AutoModelForCausalLM.from_pretrained("banana-dev/GPTrillion")
```



```
tokenizer = AutoTokenizer.from_pretrained("EleutherAI/gpt-j-6B")

prompt = (
    "In a shocking finding, scientists discovered a herd of unicorns living in  

    a remote, "  

    "previously unexplored valley, in the Andes Mountains. Even more surprising  

    to the "  

    "researchers was the fact that the unicorns spoke perfect English."  

)

input_ids = tokenizer(prompt, return_tensors="pt").input_ids

gen_tokens = model.generate(
    input_ids,
    do_sample=True,
    temperature=0.9,
    max_length=100,
)

gen_text = tokenizer.batch_decode(gen_tokens)[0]
```