

# A Comparative Study of GPT-2 and GPT-2 Based on Enhanced Self-Attention Mechanism

Wei-Hung Tu<sup>1</sup>, Neil Yen<sup>2</sup>, and Yan Pei<sup>2</sup>

<sup>1</sup> Graduate School of Computer Science and Engineering, University of Aizu,  
Aizuwakamatsu, Fukushima, 965-8580, Japan  
m5252117@u-aizu.ac.jp

<sup>2</sup> Computer Science Division, University of Aizu,  
Aizuwakamatsu, Fukushima, 965-8580, Japan  
neilyyen@u-aizu.ac.jp, peiyan@u-aizu.ac.jp

**Abstract.** In natural language processing, the quality of language models impacts applications such as machine translation and speech recognition. GPT-2, a powerful auto-regressive model with 150 million parameters, performs exceptionally well in various tasks but struggles with computational efficiency for long sequences. We have developed an optimization strategy to mitigate this issue by randomly shortening the auto-regressive length during generation. Our strategy was tested on the GPT-2 medium model using BLEU as the evaluation metric. The results revealed significant improvements in the BLEU scores, with the optimized model outperforming the original. Furthermore, the optimization also improved scores in both the top and bottom 10% of the data. Despite the promising results, there is still room for further exploration and improvement. We are currently investigating adaptive adjustments to the auto-regressive length and applying this strategy to other models, such as GPT-3. In summary, our research proposes a new strategy that enhances GPT-2's efficiency and boosts its performance, as evidenced by the improved BLEU scores. This strategy provides valuable insights for future language model optimization, holding the potential to advance the field of NLP.

**Keywords:** GPT-2, Auto regressive, Generative model, Artificial intelligence, Comparative study.

## 1 Introduction

### 1.1 Introduction to GPT-2

GPT-2, developed by OpenAI in 2019, has made a significant impact in the world of language models. It focuses on generating text that sounds natural by predicting the next word in a sequence. GPT-2 is an upgraded version of GPT, featuring more parameters, a larger model size, and better performance. The primary foundation of GPT-2 is the transformer architecture, which is a successful deep-learning model in natural language processing. One of its notable

features is the self-attention mechanism, which enables the model to capture long-distance dependencies in sequences. This capability is crucial for language modeling tasks, as words often depend on each other even when they are far apart. As stated in [1], we propose a new simple network architecture, the transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely.

One of the game-changing aspects of GPT-2 is its use of unsupervised learning for pre-training. During this stage, the model learns language patterns from massive amounts of text data without any manual labeling. As mentioned in [2], language models begin to learn these tasks without any explicit supervision when trained on a new dataset of millions of webpages called WebText. The pre-trained model can then be fine-tuned for specific NLP tasks, such as text generation, sentiment analysis, Q&A, and more. This approach fully leverages the abundance of unlabeled data available and reduces reliance on expensive human-labeled data.

GPT-2 has achieved remarkable success in various natural language processing tasks, surpassing many of the best models available at the time. With its impressive generation capabilities and versatility, it has become a go-to tool for researchers and developers addressing NLP challenges. However, it has also ignited debates surrounding potential risks, such as generating fake news and producing manipulative or unethical content.

Despite its impressive accomplishments, GPT-2 is not without flaws. There are instances where it can generate text that is grammatically correct but lacks coherence or logical consistency, and its support for different languages is not uniformly robust. Additionally, its resource-intensive nature can hinder deployment and scalability. Nevertheless, GPT-2 represents a significant milestone in the field of natural language processing, establishing a foundation for future research and applications. As stated in [2], The capacity of the language model is essential to the success of zero-shot task transfer, and increasing it improves performance in a log-linear fashion across tasks.

Since the introduction of GPT-2, OpenAI has released even more advanced models like GPT-3, which have pushed the boundaries of generation capabilities and versatility even further. However, GPT-2 remains a significant milestone in the history of NLP and continues to hold considerable research value. By improving and optimizing GPT-2, researchers can further explore the potential of self-attention mechanisms and unsupervised learning methods. As stated in [1], experiments on two machine translation tasks show these models to be superior in quality while being more parallelization and requiring significantly less time to train.

Furthermore, researchers can delve into the explainability and robustness of GPT-2 to gain a better understanding of its inner workings and make it more resilient against malicious attacks and manipulations. The success of GPT-2 in the field of natural language processing provides valuable insights and references for the development of subsequent models and technologies.

In summary, GPT-2, as an innovative and large-scale language model, has achieved significant results in the field of natural language processing and has laid a solid foundation for further research. Despite its limitations, GPT-2 continues to hold substantial research value. By improving and optimizing this model, researchers can continue to explore the potential of self-attention mechanisms and unsupervised learning methods. Understanding the background and overview of GPT-2 provides a better context for the comparative research we will discuss later, which focuses on enhancing the self-attention mechanism.

## 1.2 Research objective and motivation

In recent years, the field of natural language processing has made significant progress, largely driven by advancements in deep learning technologies. One standout model in this field is GPT-2, which utilizes self-attention and auto-regressive mechanisms to excel in various NLP tasks. However, despite its impressive performance, there is still room for improvement, particularly in terms of computational efficiency and capturing long-distance dependencies. Thus, the objectives and motivations of our research are as follows.

- Optimizing the self-attention mechanism: Our goal is to optimize the self-attention mechanism to improve computational efficiency and enhance the model’s performance in NLP tasks.
- Improving the auto-regressive mechanism: We seek to optimize the auto-regressive mechanism to address error accumulation issues and enhance the grammatical and semantic consistency of the generated text.
- Balancing local and global context information: Striking a balance between local and global context information is crucial in text generation. Overemphasizing local information can result in the loss of global consistency, while excessive attention to global information can lead to overlooked details. We will investigate techniques to achieve a balance between local and global context information within the self-attention and auto-regressive mechanisms.
- Validating the effectiveness of the improvements: Through experiments across different NLP tasks, we will compare the performance of the optimized model with the baseline model to demonstrate the effectiveness and scalability of our proposed improvements.

Our motivation lies in the potential to enhance the performance of the GPT-2 model in NLP tasks while optimizing computational efficiency. The outcomes of our research could significantly contribute to the advancement of the NLP field and offer valuable insights to researchers in related areas. Moreover, the improved model could find practical applications in intelligent dialogues, automatic summarization, knowledge graph construction, and more, providing users with higher-quality NLP services.

Throughout our research, we will employ various optimization strategies and techniques to enhance the self-attention and auto-regressive mechanisms and integrate these improvements into the GPT-2 model. Our experiments will encompass multiple NLP tasks to comprehensively evaluate the performance of

the optimized model. By analyzing the experimental results, we will assess the effectiveness and scalability of our improvements.

Finally, we will summarize our achievements, highlight the contributions we have made in optimizing the self-attention and auto-regressive mechanisms in the GPT-2 model, and discuss future research directions and potential applications. Through our research, we aspire to push the development of NLP technologies further, providing more powerful and efficient NLP solutions for real-world applications.

## 2 Experimental Design and Evaluation Methods

### 2.1 Implementation details of the optimized model

In this study, our focus is on optimizing the auto-regressive mechanism within the self-attention mechanism of the GPT-2 model. Our specific strategy involves randomly shortening certain auto-regressive lengths during the text generation process. For example, instead of considering all previous 299 words when generating the 300th word, our optimization strategy would only take into account the previous 200 words. The goal of this strategy is to reduce the model’s heavy reliance on past contexts, thereby increasing the diversity of generated text and improving generation speed.

The motivation behind this method can be summarized into three main points. Firstly, shortening the auto-regressive length helps to promote diversity in the generated text and prevents the production of excessively repetitive content. Secondly, it reduces computational load, resulting in faster text generation. Lastly, this strategy enhances the model’s ability to generalize, enabling it to adapt better to new samples and text styles that may differ from the training data.

In our forthcoming experiments, we will explore the impact of this strategy on the generative capabilities of the GPT-2 model. We will conduct a series of experiments to analyze how different auto-regressive lengths affect the quality of generated content. Additionally, we will evaluate the effectiveness of this method across various natural language processing tasks. Through these experiments, we aim to provide valuable insights into language model optimization and offer guidance for future improvements.

In conclusion, this research focuses on optimizing the auto-regressive mechanism within the self-attention mechanism of the GPT-2 model by randomly shortening auto-regressive lengths during text generation. This approach has the potential to enhance the diversity of generated text, accelerate the generation process, and improve the model’s generalization capability. Subsequent experiments will validate the effectiveness and feasibility of this strategy, providing empirical support for further optimization of language models.

### 2.2 Evaluation metrics and experimental setup

We use BLEU (Bilingual Evaluation Understudy) as our primary evaluation metric in this study. BLEU is commonly used to measure the grammatical and

semantic similarity between generated text and manually written reference text. It relies on n-gram precision for scoring. The BLEU score ranges from 0 to 1, with higher values indicating a higher similarity between the model-generated text and the human-written reference text, thus reflecting better model performance.

For our experiments, we selected the GPT-2 medium model configuration, which consists of 1.5 billion parameters. Among different versions of GPT-2, such as GPT-2 small or GPT-2 large, the medium version strikes a suitable balance between computational requirements and performance. It exhibits strong performance across various natural language processing tasks, making it an appropriate baseline model for optimizing the auto-regressive mechanism within the self-attention mechanism.

To accurately evaluate the impact of our optimization, we will compare the optimized GPT-2 medium model with the original GPT-2 medium model. The original model serves as our baseline, and we expect the optimized model to outperform it in terms of BLEU scores.

This experimental setup allows us to precisely measure the effect of optimizing the auto-regressive mechanism within the self-attention mechanism on the generative capabilities of the GPT-2 medium model. By adopting this approach, we aim to enhance the credibility of our results and contribute to the advancements in the field of natural language processing. We are excited about the potential insights and possibilities this research may bring to the self-attention mechanism and its application in auto-regressive models.

### 3 Experimental Results and Analysis

In terms of BLEU scoring, the optimized GPT-2 model demonstrated significant improvements compared to the original GPT-2 model. The average BLEU score of our optimized model was 2.902508133906813e-06, while the original GPT-2 model had an average score of only 2.803913135029713e-158 (Table 1). This result confirms the effectiveness of our optimization strategy and highlights the superior generative capabilities of our optimized model.

Further analysis of our data revealed that in the top 10% of the generated text in Table 2, the average BLEU score of the optimized model reached 2.9156350551154368e-05, which is significantly higher than the original GPT-2 model's score of 2.81659414920824e-157. This finding indicates that our optimization strategy successfully enhanced the GPT-2 model's ability to generate high-quality text in the top percentile.

However, in the lowest 10% of the generated text, our optimized model's performance was similar to that of the original model, with BLEU scores of 4.797929067427243e-232 and 4.808473266007582e-232, respectively. This suggests that our optimization strategy has a limited impact on improving the lower-end generative capabilities of the model. We hypothesize that this may be because our optimization strategy primarily focuses on enhancing the model's top-end performance while having a limited effect on the lower end.

**Table 1.** BLEU scores of optimized and original GPT-2 models

Metric	Average
Optimized GPT-2	$2.902508133906813 \times 10^{-6}$
Original GPT-2	$2.803913135029713 \times 10^{-158}$

**Table 2.** Average BLEU scores (top 10% and bottom 10%) of optimized and original GPT-2 models

Metric	Top 10% Average	Bottom 10% Average
Optimized GPT-2	$2.9156350551154368 \times 10^{-5}$	$4.797929067427243 \times 10^{-232}$
Original GPT-2	$2.81659414920824 \times 10^{-157}$	$4.808473266007582 \times 10^{-232}$

## 4 Conclusions and Future Work

### 4.1 Summary of research results

Upon reviewing our experimental results, we can conclude that our optimization strategy has indeed improved the generative capabilities of the GPT-2 model to a certain degree. This improvement is particularly notable in the top 10% of the generated text, where our optimized model achieved a significantly higher BLEU score compared to the original model. This outcome provides strong evidence for the effectiveness of our optimization strategy.

However, it is important to acknowledge that our approach showed limited success in enhancing the lower-end performance of the model. This finding highlights the need for further refinement and exploration in our optimization strategy to address this aspect of the model’s generative capabilities.

These results emphasize the importance of achieving a balanced performance across both high-quality and low-quality segments of generated text in future optimization efforts. It is crucial to strive for comprehensive robustness in the model’s performance across all aspects of text generation, rather than focusing solely on improving the high-quality domain. By addressing the limitations identified in our study, we can pave the way for more robust and reliable language models in the future.

### 4.2 Future work

In our future work, our primary focus will be on refining our optimization strategy, specifically targeting the improvement of the model’s generative capabilities at the lower-end performance range. We will explore various methods for adjusting the self-attention and autoregressive mechanisms to enhance the model’s proficiency in generating lower-quality text segments, while still maintaining the quality of high-end outputs.

Additionally, we aim to incorporate additional evaluation metrics, such as ROUGE and METEOR, to provide a more comprehensive assessment of the

performance of our optimized model. By considering multiple evaluation metrics, we can gain a more nuanced understanding of the model's strengths and weaknesses.

Moreover, we have plans to extend our optimization strategy to other generative models, such as GPT-3 and GPT-4, to examine the applicability and effectiveness of our approach across a broader range of models.

In conclusion, our research has successfully optimized the GPT-2 model, offering new perspectives for future investigations. We are excited to continue improving the generative capabilities of the GPT-2 model in our upcoming work and to expand the application of our findings to a wide range of scenarios and domains.

## References

1. A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
2. A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.