

**Automated Title Generation for Research Papers: A Transformer-Based  
Approach**

by

**Maryam Hosseinzadeh**

A Project Submitted in Partial Fulfillment of the Requirements for the Degree of

Master of Engineering

In The Department of Electrical and Computer Engineering

© Maryam Hosseinzadeh, 2024  
University of Victoria

All rights reserved. This project may not be reproduced in whole or in part, by photocopying or other means, without the permission of the author.

Supervisory Committee

---

Dr. Amirali Baniyadi, **Supervisor**

(Department of Electrical and Computer Engineering)

---

Dr. Fayez Gebali, **Committee Member**

(Department of Electrical and Computer Engineering)

## ABSTRACT

The task of generating concise, informative, and relevant titles for research papers is a critical but challenging aspect of academic publishing. This project report explores the application of Transformer-based Large Language Models (LLMs) to automatically generate and suggest research paper titles from their abstracts. Building on the foundational Transformer architecture introduced by Vaswani et al., we evaluate the performance of different models, including a custom-built LLM, a pre-trained GPT-2 model, and a fine-tuned version of GPT-2. Through qualitative analysis, we demonstrate that fine-tuning GPT-2 on a specific dataset of research paper abstracts and titles significantly enhances the coherence, relevance, and contextual accuracy of the generated titles. We address challenges such as hallucinations in LLM-generated text and discuss the importance of high-quality datasets and task-specific fine-tuning. This work contributes to the broader understanding of the capabilities and limitations of LLMs in specialized NLP tasks, offering insights for future research and applications in academic publishing.

# Contents

<b>Abstract</b>	<b>iii</b>
<b>Table of Contents</b>	<b>iv</b>
<b>List of Tables</b>	<b>vii</b>
<b>List of Figures</b>	<b>viii</b>
<b>Acronyms</b>	<b>ix</b>
<b>Acknowledgements</b>	<b>x</b>
<b>Dedication</b>	<b>xi</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Problem Statement . . . . .	1
1.2 Overview . . . . .	1
1.3 Motivation . . . . .	2
1.4 Contribution . . . . .	2
1.5 project Structure . . . . .	2
<b>2 Background</b>	<b>4</b>
2.1 Introduction to NLP . . . . .	4
2.2 Sequence Transduction Models . . . . .	4
2.2.1 Introduction to RNN . . . . .	4
2.2.2 CNNs . . . . .	5
2.3 The Transformer Architecture . . . . .	5
2.4 Introduction to LLMs . . . . .	5
2.4.1 Generative Pre-trained Transformer 2 (GPT-2) . . . . .	5
2.4.2 Bidirectional Encoder Representations from Transformers (BERT)	5

2.5	Title Generation in Academic Publishing . . . . .	6
2.6	Previous Work . . . . .	7
<b>3</b>	<b>Methodology</b>	<b>8</b>
3.1	Introduction to Tools and Environment . . . . .	8
3.1.1	Computational Resources . . . . .	8
3.1.2	Programming Language and Libraries . . . . .	9
3.2	Introduction to Transformers . . . . .	9
3.2.1	Positional Encodings . . . . .	9
3.2.2	Attention Mechanism . . . . .	10
3.2.3	Self-Attention . . . . .	10
3.3	Overview of the Transformer Model Architecture . . . . .	11
3.4	Approach . . . . .	12
3.5	Dataset . . . . .	12
3.5.1	Considerations . . . . .	14
3.6	Implementing the Attention Networks and LLM from Scratch . . . . .	15
3.6.1	Dataset Preparation . . . . .	15
3.6.2	Tokenization . . . . .	15
3.6.3	Model Architecture . . . . .	15
3.6.4	Building the Transformer . . . . .	16
3.6.5	Training the Model . . . . .	16
3.6.6	Implementation Challenges and Solutions . . . . .	17
3.7	Transfer Learning with GPT-2 . . . . .	17
3.7.1	Transfer Learning . . . . .	17
3.7.2	Why Transfer Learning is a Good Choice . . . . .	18
3.7.3	GPT-2: An Overview . . . . .	18
3.7.4	Application of GPT-2 in Our Project . . . . .	19
3.7.5	Results and Conclusion . . . . .	20
<b>4</b>	<b>Discussion and Results</b>	<b>21</b>
4.1	Effectiveness of LLMs . . . . .	21
4.2	LLM Hallucinations and Other Factors . . . . .	21
4.3	Results . . . . .	22
4.4	Baseline LLM Performance . . . . .	22
4.5	GPT-2 Performance . . . . .	22

4.6	Fine-Tuned GPT-2 Performance . . . . .	23
4.7	Case Study: Generated Titles . . . . .	23
4.8	Conclusion . . . . .	23
<b>5</b>	<b>Conclusion</b>	<b>26</b>
5.1	Summary of Findings . . . . .	26
5.2	Effectiveness of the Transformer Architecture . . . . .	26
5.3	Addressing LLM Hallucinations and Challenges . . . . .	27
5.4	Implications for Future Research . . . . .	27
5.5	Final Remarks . . . . .	28
	<b>Bibliography</b>	<b>29</b>

# List of Tables

Table 3.1 Specifications of the Research Paper Dataset 2023 . . . . .	13
Table 4.1 Comparison of Titles Generated by Different Models . . . . .	25

# List of Figures

Figure 2.1 The Transformer architecture [1]. The model consists of an encoder and a decoder, each composed of a stack of layers. . . . .	6
--	---



## Acronyms

- **AI**: Artificial Intelligence
- **ML**: Machine Learning
- **RNN**: Recurrent Neural Network
- **CNN**: Convolutional Neural Network
- **GPT**: Generative Pre-trained Transformer
- **NLP**: Natural Language Processing
- **LLM**: Large Language Model
- **BERT**: Bidirectional Encoder Representations from Transformers

## ACKNOWLEDGEMENTS

First and foremost, I want to thank my dearest **supervisor, Professor Amirali Baniyadi**, for his unwavering support, belief in me, and for giving me this incredible opportunity. His guidance has profoundly impacted my life, and I am deeply grateful for everything he has done for me during my time at the University of Victoria. Thank you, Professor Baniyadi, for your consistent support and encouragement.

Next, I want to extend my heartfelt gratitude to **my parents** and **my sister** for always being there for me and for all the sacrifices they have made to allow me to pursue my passion and dreams. I could not have done it without you!

I am also grateful to **Javad** and **Nasrin** for their extraordinary and unexpected generosity, as well as for their efforts to advise me through both academic pursuits and difficult life choices.

I would like to extend my heartfelt gratitude to **Nikoo** and **Naeeme**. They have been a true lifesaver, offering sisterly support, and my appreciation for their assistance will last a lifetime. Even if I expressed my thanks daily, it would never suffice!

I want to thank **Amirhossein**, who was my first trustworthy person in this city and helped me to settle in this new environment and at the university.

Finally, I want to express my deepest gratitude to **Ardeshir**. He has guided, supported, and taught me in exceptional ways. His consideration throughout this journey has made it much more enjoyable and has helped me discover new strengths within myself to grow and develop.

I would like to thank **Golnaz**, **Negar** and **Parnian**, along with all my friends who have been great companions and supported me through my life journey. And to everyone who has helped me along the way—thank you!

## DEDICATION

This project is dedicated to my parents and my sister, whose unconditional love, support, and sacrifices have made this journey possible. Their unwavering belief in me and their constant encouragement have been my greatest source of strength. I am forever grateful for their endless patience, understanding, and for always standing by my side. This work is a testament to their immense contributions to my life and my academic pursuits.

# Chapter 1

## Introduction

### 1.1 Problem Statement

Generating titles for research papers is a critical yet challenging task in the academic publishing process. Titles need to be concise, informative, and reflective of the paper's content to attract readership and convey the essence of the research. Traditional methods of title generation often rely on manual effort, which can be time-consuming and inconsistent. The advent of LLMs and the Transformer architecture presents an opportunity to automate this task with high accuracy and efficiency. This project aims to explore the effectiveness of LLMs, particularly the Transformer-based models, in automating the generation of research paper titles based on their abstracts [2, 3].

### 1.2 Overview

This project investigates the application of Transformer-based LLMs for the task of generating research paper titles. We evaluate different model configurations, including a custom-built LLM, the pre-trained GPT-2 model, and a fine-tuned version of GPT-2 on a specific dataset of research paper abstracts and titles. The study involves qualitative analysis of the generated titles to assess their coherence, relevance, and contextual accuracy [4].

## 1.3 Motivation

The motivation behind this research stems from the need to streamline the academic publishing process and enhance the quality of research paper titles. Manual title generation is prone to variability and inefficiency, which can be addressed through automation. LLMs, with their ability to understand and generate human-like text, offer a promising solution. By leveraging advanced models like the Transformer, we aim to improve the consistency and accuracy of title generation, thereby benefiting researchers and publishers alike [2].

## 1.4 Contribution

This project makes several key contributions to the field of Natural Language Processing (NLP) and academic publishing:

- Demonstrates the application of Transformer-based LLMs for the task of generating research paper titles.
- Provides a qualitative analysis of the effectiveness of different LLM configurations, highlighting the improvements achieved through fine-tuning.
- Identifies and addresses challenges such as hallucinations in LLM-generated text, offering insights for future research.
- Contributes to the broader understanding of the capabilities and limitations of LLMs in specialized NLP tasks [3, 4, 2].

## 1.5 project Structure

The subsequent chapters of this project unfold a comprehensive examination of the methodology, results, and implications of our study on the effectiveness of Transformer-based large language models (LLMs) for generating research paper titles from abstracts.

Chapter 2 provides the necessary background information, setting the stage for the detailed discussions to follow. This chapter explores the historical development and foundational concepts of sequence transduction models, with a particular focus on the evolution from recurrent neural networks (RNNs) and convolutional neural networks

(CNNs) to the revolutionary Transformer architecture. It also delves into the core principles of self-attention mechanisms, which form the backbone of the Transformer model, enabling it to capture complex dependencies within sequential data efficiently.

In Chapter 3, we delve into the methodology employed in our study. This chapter begins with an in-depth introduction to the Transformer architecture and the attention mechanism, highlighting their significance in modern NLP tasks. We describe our specific approach to training and fine-tuning the models on a carefully curated dataset of research paper abstracts and titles. This includes detailed explanations of the preprocessing steps, model training procedures, and evaluation techniques used to ensure robust and reliable results. The chapter also discusses the various configurations of the models, including a custom-built LLM, the pre-trained GPT-2 model, and the fine-tuned GPT-2 model, providing insights into their respective training processes and performance metrics.

Chapter 4 presents a qualitative discussion of the results obtained from our experiments. Through detailed comparisons, we illustrate the performance differences between the custom-built LLM, the pre-trained GPT-2, and the fine-tuned GPT-2 models. Examples are provided to demonstrate how the fine-tuned GPT-2 model generates more coherent and contextually accurate titles compared to the other models. This chapter also addresses several challenges encountered during the study, such as the phenomenon of hallucinations where models generate plausible but incorrect outputs. The importance of high-quality datasets and the critical role of model fine-tuning in mitigating these issues are emphasized, offering valuable lessons for future research.

Finally, Chapter 5 concludes the project by summarizing the key findings and discussing their broader implications. This chapter highlights the effectiveness of Transformer-based LLMs in the specific task of title generation for research papers, acknowledging the challenges faced and the solutions implemented. It reflects on the potential for further advancements in this area, suggesting avenues for future research that could build on the foundation laid by this study. The chapter concludes with final remarks on the transformative impact of our work, underscoring the contributions made to the field of NLP and the potential for ongoing innovation.

# Chapter 2

## Background

### 2.1 Introduction to NLP

NLP is a subfield of artificial intelligence (AI) that focuses on the interaction between computers and humans through natural language. The goal of NLP is to enable computers to understand, interpret, and generate human language in a way that is both meaningful and useful. Key applications of NLP include machine translation, sentiment analysis, text summarization, and question answering [2].

### 2.2 Sequence Transduction Models

Sequence transduction models are a class of models used in NLP to convert sequences of one type into sequences of another type. Examples include translating sentences from one language to another, converting speech signals to text, and summarizing long documents. Traditionally, sequence transduction tasks have been tackled using recurrent neural networks (RNNs) and CNNs [4].

#### 2.2.1 Introduction to RNN

RNNs are a type of neural network designed for sequential data. They process input sequences one element at a time, maintaining a hidden state that captures information about previous elements. Despite their ability to handle sequences of varying lengths, RNNs suffer from limitations such as difficulty in capturing long-range dependencies due to vanishing gradient problems [2].

### 2.2.2 CNNs

CNNs, originally developed for image processing tasks, have also been applied to sequence transduction. CNNs use convolutional layers to capture local dependencies in the input sequences. Although CNNs can be parallelized more easily than RNNs, they often require a large number of layers to capture long-range dependencies effectively [4].

## 2.3 The Transformer Architecture

The Transformer architecture, introduced by Vaswani et al. in the paper "Attention is All You Need" [1], represents a significant departure from traditional sequence transduction models. Unlike RNNs and CNNs, Transformers rely entirely on self-attention mechanisms to process input sequences, as illustrated in Figure 2.1.

## 2.4 Introduction to LLMs

LLMs are a subset of NLP models that are trained on massive datasets to understand and generate human language. Examples of LLMs include GPT-2, BERT, and T5. These models leverage the Transformer architecture to perform a wide range of NLP tasks with high accuracy [2].

### 2.4.1 Generative Pre-trained Transformer 2 (GPT-2)

GPT-2, developed by OpenAI, is a large-scale Transformer-based language model. It is trained on diverse internet text data to generate coherent and contextually relevant text. GPT-2 can be fine-tuned on specific datasets to adapt to specialized tasks, such as generating titles for research papers [3].

### 2.4.2 Bidirectional Encoder Representations from Transformers (BERT)

BERT, introduced by Google, is a Transformer-based model designed for a variety of NLP tasks. Unlike GPT-2, which is primarily generative, BERT is designed to understand the context of words in a bidirectional manner, making it highly effective for tasks such as text classification and question answering [4].



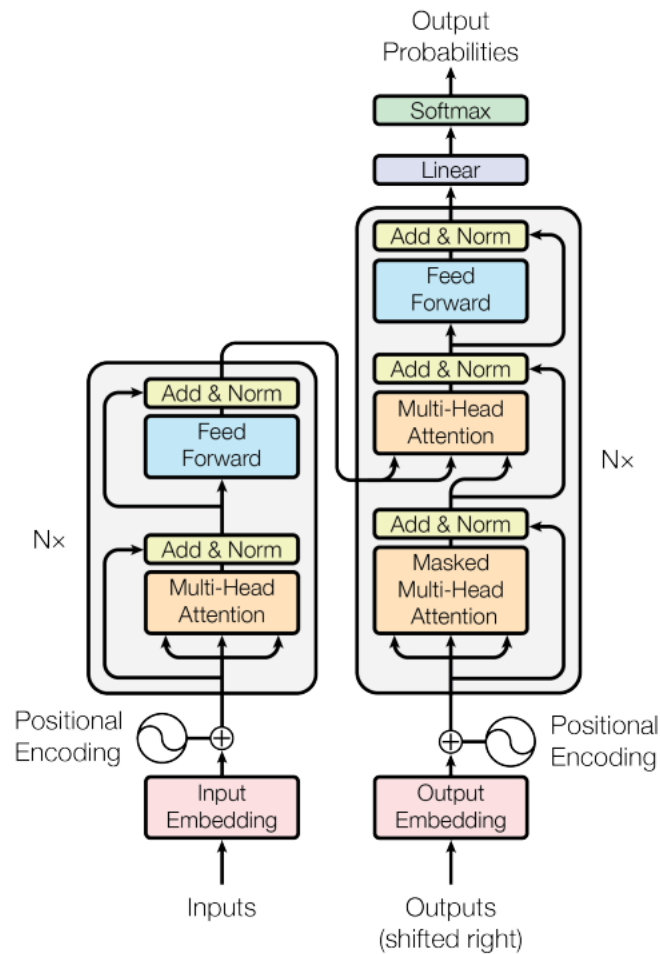


Figure 2.1: The Transformer architecture [1]. The model consists of an encoder and a decoder, each composed of a stack of layers.

## 2.5 Title Generation in Academic Publishing

Title generation for research papers is a crucial aspect of academic publishing. A well-crafted title can attract readers, convey the essence of the research, and improve the visibility of the paper. Automating this task using LLMs can significantly enhance the efficiency and consistency of the process [2].

## 2.6 Previous Work

Several studies have explored the use of LLMs and Transformer architectures for various NLP tasks. Research has demonstrated the effectiveness of these models in generating human-like text, translating languages, summarizing documents, and more. However, the application of LLMs for title generation in academic publishing remains an emerging area of study [3, 4, 2].

# Chapter 3

## Methodology

### 3.1 Introduction to Tools and Environment

In this project, we utilized a combination of powerful computational tools to implement and train the Transformer-based model for automated title generation. The primary tools and environment details are outlined below.

#### 3.1.1 Computational Resources

The core computations for this project were performed on a high-performance workstation equipped with advanced hardware capable of handling the intensive demands of deep learning tasks. The system was configured with the following specifications:

- **GPU:** NVIDIA GeForce RTX 3090, offering substantial computational power for accelerated model training, particularly for large-scale transformer models.
- **CPU:** AMD Ryzen 9 with 24 cores, providing robust parallel processing capabilities to efficiently manage data preprocessing and model evaluation tasks.
- **Memory:** 128 GB of RAM, ensuring smooth handling of large datasets and complex model architectures without memory bottlenecks.

All data processing, model training, and evaluation were conducted on this workstation, ensuring consistent performance and minimal latency. The powerful GPU was particularly beneficial for managing the computationally intensive operations required for training transformer models, while the multi-core CPU facilitated efficient

handling of concurrent tasks. The ample memory capacity allowed for the processing of large datasets without encountering memory limitations.

### 3.1.2 Programming Language and Libraries

Python was the primary programming language used for implementing the models and processing the data. Python’s rich ecosystem of libraries provided the necessary tools for developing and fine-tuning the Transformer models.

The Python environment was set up on the computer, leveraging its computational power to execute the intensive training processes involved in fine-tuning the Transformer model.

## 3.2 Introduction to Transformers

Transformers have fundamentally transformed the field of natural language processing (NLP) since their introduction by Vaswani et al. in the seminal paper “*Attention is All You Need*” [1]. The model’s unprecedented success across various NLP tasks, including translation, summarization, and title generation, can be attributed to three key innovations: positional encodings, attention mechanisms, and a specialized variant of attention known as self-attention. These innovations enable transformers to effectively model complex dependencies in sequential data by dynamically assigning different levels of importance to words in a sequence.

### 3.2.1 Positional Encodings

In traditional models like Recurrent Neural Networks (RNNs), the sequential nature of data processing inherently preserved word order. However, transformers, which process words in parallel, needed a way to understand the order of words within a sequence. This challenge is elegantly addressed by *positional encodings*.

Positional encodings involve assigning a unique positional index to each word in a sentence before it is fed into the neural network. For example, when translating text from English to French, each word is tagged with a positional identifier—1, 2, 3, and so on—reflecting its position within the sentence. This positional information is then embedded into the input data rather than being managed by the network’s architecture.

During training, the transformer model learns to interpret these positional encodings, understanding the significance of word order directly from the data. This innovation not only preserves the sequential nature of language but also enhances the model’s ability to capture contextual dependencies, making transformers more efficient and easier to train compared to RNNs.

### 3.2.2 Attention Mechanism

The second innovation, *attention*, has become a cornerstone in modern machine learning, particularly in tasks involving translation and sequence-to-sequence learning.

Consider the sentence: “The agreement on the European Economic Area was signed in August 1992.” Translating this sentence into French is not as straightforward as translating each word individually. In French, word order may change, and grammatical rules such as gender agreement must be respected—‘European’ might precede ‘economic,’ and the word ‘européenne’ must be in its feminine form to agree with ‘la zone.’

The attention mechanism allows the model to focus on relevant parts of the input sentence when generating each word of the output sentence. Instead of translating each word in isolation, the model evaluates the entire sentence to determine which words are most relevant to the current translation step. This dynamic focus on different parts of the sentence enables the model to learn complex linguistic features, such as gender, word order, and plurality, through exposure to large amounts of parallel data.

### 3.2.3 Self-Attention

While attention mechanisms had been introduced prior to transformers, the true innovation of the transformer model lies in its application of *self-attention*. This approach extends the concept of attention by enabling the model to relate different positions within the same input sentence to each other, thereby understanding the contextual meaning of words.

Self-attention empowers the model to build rich internal representations of language. For instance, the model might learn that the terms “programmer,” “software engineer,” and “software developer” are often used interchangeably. It might also grasp the rules of grammar, gender, and tense, refining its understanding as it processes vast amounts of text data.

Consider the sentences: “Server, can I have the check?” and “Looks like I just crashed the server.” The word “server” carries different meanings in these contexts. Through self-attention, the model can discern the correct interpretation by focusing on the surrounding words—‘check’ in the first sentence and ‘crashed’ in the second. This ability to disambiguate meaning based on context is one of the many strengths of self-attention.

In summary, transformer models are built on the innovations of positional encodings, attention, and self-attention. These mechanisms work together to enable the model to understand and generate language with remarkable accuracy and efficiency, laying the groundwork for a wide range of NLP applications.

### 3.3 Overview of the Transformer Model Architecture

The transformer model consists of an encoder and a decoder. The encoder is composed of a stack of identical layers, each containing two main sub-layers: multi-head self-attention and a feedforward network. The decoder has an additional third sub-layer that performs multi-head attention over the encoder’s output.

#### Encoder

The encoder is composed of a stack of  $N$  identical layers. Each layer has two sub-layers:

- **Multi-Head Self-Attention Mechanism:** This allows the encoder to attend to different positions of the input sequence simultaneously.
- **Position-wise Fully Connected Feed-Forward Network:** This consists of two linear transformations with a ReLU activation in between.

Residual connections are employed around each of the two sub-layers, followed by layer normalization. The output of each sub-layer is  $\text{LayerNorm}(x + \text{Sublayer}(x))$ , where  $\text{Sublayer}(x)$  is the function implemented by the sub-layer itself.

#### Decoder

The decoder is also composed of a stack of  $N$  identical layers. In addition to the two sub-layers in each encoder layer, the decoder inserts a third sub-layer, which performs

multi-head attention over the output of the encoder stack. Similar to the encoder, residual connections are employed around each of the sub-layers, followed by layer normalization.

To prevent positions from attending to subsequent positions, the self-attention sub-layers in the decoder mask future positions (setting them to  $-\infty$  before the softmax step). This masking, combined with the fact that the output embeddings are offset by one position, ensures that the predictions for position  $i$  can depend only on the known outputs at positions less than  $i$ .

### 3.4 Approach

Our approach involves training a transformer-based model to generate titles for research papers based on their abstracts. The model architecture follows the standard transformer setup with modifications to suit the title generation task. The overall process is as follows:

1. **Data Preprocessing:** The dataset containing research paper titles and abstracts is tokenized using a Byte Pair Encoding (BPE) tokenizer. This tokenizer splits words into subword units, enabling the model to handle rare words effectively.
2. **Model Training:** The transformer model is trained on the tokenized dataset. The training process involves minimizing the cross-entropy loss between the predicted and actual titles. The optimizer used is Adam with a learning rate scheduler that implements a warm-up strategy followed by a decay phase.
3. **Evaluation:** The model's performance was evaluated based on the coherence, relevance, and contextual accuracy of the generated titles.

### 3.5 Dataset

The dataset utilized in this project is titled "Research Paper Dataset 2023," comprising details pertinent to research papers, including titles and abstracts. The dataset includes the following attributes:

- **Title (dtype: string):** The title of the research paper.

- **Abstract (dtype: string):** The abstract of the research paper.

The dataset is organized into a single split, detailed below:

Dataset Split	Details
<b>Train Split</b>	
Name	train
Size (Bytes)	2,363,569,633
Number of Examples	2,311,491
<b>Download Information</b>	
Download Size	1,423,881,564 bytes
Total Dataset Size	2,363,569,633 bytes

Table 3.1: Specifications of the Research Paper Dataset 2023

This dataset is available for public access and can be downloaded from the following link: [https://huggingface.co/datasets/Falah/research\\_paper2023](https://huggingface.co/datasets/Falah/research_paper2023). It is provided under the Apache License 2.0 [5].

**Applications for NLP Text Classification and Chatbot Models:** The "Research Paper Dataset" serves as a valuable resource for a variety of NLP tasks, such as text classification and generating book titles within chatbot models. Here are some ways to utilize this dataset:

**Text Classification:** By leveraging the titles and abstracts, this dataset can be employed to train a model that classifies text. This involves labeling research papers based on their subject areas or fields of study, enabling the model to categorize new research papers accurately. For instance, the model could predict if a paper pertains to fields like computer science, biology, or physics. This trained model can also be adapted for other text categorization tasks.

**Generating Book Titles for Chatbots:** Using the titles from research papers, a natural language generation model (such as a sequence-to-sequence model or a transformer-based model) can be trained to create book titles. Fine-tuning the model on these titles helps it learn the patterns and structures necessary for generating relevant and meaningful titles. This functionality can enhance chatbot models, enabling them to recommend books based on specific research topics or interests.



## Potential Benefits

- **Enhanced Recommendations:** Chatbots can offer more tailored and relevant book recommendations by generating titles related to specific research topics.
- **Improved User Interaction:** Integrating text classification models allows chatbots to better understand user queries, resulting in more accurate and engaging responses.
- **Efficient Knowledge Discovery:** Researchers and students can use the text classification model to categorize vast amounts of research papers quickly, aiding in the swift retrieval of relevant information.

### 3.5.1 Considerations

Throughout this project, several important considerations were addressed to ensure the quality, performance, and ethical implications of our work:

- **Data Preprocessing:** We carefully implemented proper preprocessing steps, including text cleaning, tokenization, and encoding, to prepare the dataset for model training. These steps were crucial to ensure that the data fed into the model was clean, structured, and suitable for learning.
- **Model Selection and Tuning:** We selected appropriate model architectures and hyperparameters that were well-suited to our specific task of generating titles from research paper abstracts. Fine-tuning these models was essential to achieve optimal performance and to adapt the pre-trained models to our dataset's unique characteristics.
- **Ethical Usage:** We made a concerted effort to ensure the responsible and ethical use of the generated titles and text classification predictions. This involved respecting copyright and intellectual property rights, and avoiding the use of generated outputs in a way that could infringe upon these rights. Additionally, we aimed to prevent the dissemination of misleading or incorrect information, adhering to high ethical standards in our research and its applications.

## 3.6 Implementing the Attention Networks and LLM from Scratch

In this section, we detail the process of building a LLM from scratch using PyTorch, specifically tailored for generating titles from research paper abstracts. This implementation draws inspiration from the foundational principles of transformer architecture as introduced by Vaswani et al. in the seminal paper "Attention is All You Need" [1].

### 3.6.1 Dataset Preparation

The first step in developing our model involved acquiring and preparing the dataset. We utilized the "Research Paper Dataset 2023," which contains research paper abstracts and their corresponding titles. This dataset, sourced from Hugging Face, was divided into training and validation splits to facilitate effective model training and evaluation.

### 3.6.2 Tokenization

Transformers require numerical input, necessitating the conversion of text data into numerical representations through tokenization. We employed a Byte Pair Encoding (BPE) tokenizer to achieve this. The tokenizer was trained on the corpus data to generate a vocabulary specific to the research paper abstracts and titles. This process involves:

- **Training the Tokenizer:** Separate tokenizers were trained for the source (abstracts) and target (titles) languages.
- **Generating Vocabulary:** The tokenizer creates a vocabulary of unique tokens, essential for converting text into token IDs that the model can process.

### 3.6.3 Model Architecture

The core of our implementation is the transformer model, which leverages self-attention mechanisms to process input sequences. The architecture consists of the following components:

- **Embedding Layer:** Converts token IDs into dense vectors that capture the semantic meaning of the tokens.
- **Positional Encoding:** Adds information about the position of tokens within the sequence to the embeddings, addressing the lack of order awareness in parallel processing.
- **Multi-Head Attention:** Facilitates the model's ability to focus on different parts of the input sequence, enhancing its understanding of contextual relationships.
- **Feedforward Network:** Applies non-linear transformations to the data, allowing the model to learn complex patterns.
- **Layer Normalization and Residual Connections:** Stabilize and improve the training process by normalizing inputs and enabling gradient flow through skip connections.

### 3.6.4 Building the Transformer

To build the transformer, we assembled the individual components described above. The process involves creating the encoder and decoder blocks, each consisting of multiple layers of multi-head attention and feedforward networks. The encoder processes the input abstracts, while the decoder generates titles based on the encoder's output and the previously generated tokens.

### 3.6.5 Training the Model

Training the model involves several steps to ensure it learns to generate coherent and relevant titles from abstracts:

1. **Data Loading and Preprocessing:** The dataset is loaded and preprocessed, converting text to token IDs using the trained tokenizers.
2. **Forward Pass:** The input abstracts are passed through the encoder to generate contextual embeddings. These embeddings are then fed into the decoder, which generates the output sequence (titles) token by token.

3. **Loss Calculation:** The model's predictions are compared to the actual titles using a loss function, typically Cross-Entropy Loss, to quantify the prediction error.
4. **Backpropagation and Optimization:** The gradients of the loss with respect to the model parameters are computed and used to update the parameters, minimizing the loss over time.
5. **Validation:** After each training epoch, the model is evaluated on the validation dataset to monitor its performance and prevent overfitting.

### 3.6.6 Implementation Challenges and Solutions

During the implementation, several challenges were encountered and addressed:

- **Handling Long Sequences:** The maximum sequence length was set to ensure efficient processing while retaining relevant information.
- **Preventing Overfitting:** Techniques such as dropout and regularization were applied to improve the model's generalization.
- **Mitigating Model Hallucinations:** Strategies like fine-tuning with specific datasets and implementing reinforcement learning from human feedback (RLHF) were explored to reduce the occurrence of implausible outputs.

## 3.7 Transfer Learning with GPT-2

### 3.7.1 Transfer Learning

Transfer learning is a machine learning technique where a model developed for a specific task is reused as the starting point for a model on a second task. This method leverages the knowledge gained from solving one problem and applies it to a different but related problem. This approach is particularly effective when the second task has limited training data, as it allows the model to benefit from the extensive training data used for the first task.

The primary advantage of transfer learning is that it can significantly improve the performance of a model on the target task by utilizing pre-existing knowledge. This is especially beneficial in scenarios where data collection is challenging or expensive. By

transferring learned features from a pre-trained model, we can achieve better results with less computational resources and time.

### 3.7.2 Why Transfer Learning is a Good Choice

In the context of our project, transfer learning was an optimal choice due to several reasons:

1. **Data Scarcity:** Our dataset, while substantial, is not as large as the datasets typically required to train complex models from scratch. By leveraging a pre-trained model, we can overcome the limitations posed by our dataset size.

2. **Training Efficiency:** Training deep learning models from scratch requires significant computational power and time. Transfer learning allows us to start from a well-trained model, drastically reducing the training time and computational cost.

3. **Performance Improvement:** Pre-trained models have already learned useful features from large and diverse datasets. By fine-tuning these models on our specific task, we can enhance their performance, often achieving state-of-the-art results with minimal effort.

### 3.7.3 GPT-2: An Overview

GPT-2 (Generative Pre-trained Transformer 2) is a state-of-the-art language model developed by OpenAI. It is designed to generate human-like text by predicting the next word in a sentence given all previous words. GPT-2 is based on the transformer architecture, which has become the foundation for many recent advances in NLP.

#### How GPT-2 Works

GPT-2 operates using a transformer architecture, which relies on self-attention mechanisms to process input text. The self-attention mechanism allows the model to weigh the importance of different words in a sentence relative to each other, enabling it to capture complex dependencies and contextual information. GPT-2 is trained in an unsupervised manner on a diverse and extensive corpus of internet text, allowing it to learn a wide range of linguistic patterns and knowledge.

The architecture of GPT-2 consists of multiple layers of transformers, each layer refining the model's understanding of the input text. During training, GPT-2 learns to predict the next word in a sequence, optimizing its parameters to minimize the

difference between its predictions and the actual next words. This process endows GPT-2 with the ability to generate coherent and contextually relevant text.

### **Training Objective and Dataset Diversity**

The training of GPT-2 on such a large and diverse dataset allows it to generalize well to various language tasks. The model learns to predict the next word in a sequence by minimizing the difference between its predictions and the actual next words during training. This process, called unsupervised learning, equips GPT-2 with a robust understanding of language structure and content, enabling it to generate coherent and contextually relevant text.

The dataset used for training GPT-2 contains a vast array of information from numerous web pages, covering topics from different domains. This diversity ensures that the model is exposed to a wide range of linguistic patterns, styles, and contexts, which enhances its ability to understand and generate text across various subjects.

### **3.7.4 Application of GPT-2 in Our Project**

For our project, we utilized transfer learning by fine-tuning GPT-2 on our specific dataset. The pre-trained GPT-2 model served as a robust starting point, providing a strong foundation of linguistic knowledge. We then fine-tuned this model on our dataset to adapt it to our particular task requirements.

The steps involved in this process were as follows:

1. **Pre-training:** We began with GPT-2, which had already been pre-trained on a diverse corpus of text data. This pre-training phase equipped the model with a broad understanding of language.

2. **Fine-tuning:** We fine-tuned GPT-2 on our dataset, allowing the model to adapt its learned features to our specific task. This involved additional training on our dataset, adjusting the model's weights to optimize performance on our data.

3. **Evaluation:** After fine-tuning, we evaluated the performance of the GPT-2 model on our task. The results indicated a significant improvement in performance compared to models trained from scratch or using less sophisticated techniques.

The use of GPT-2 through transfer learning enabled us to achieve better results efficiently. The model's pre-existing knowledge, combined with the fine-tuning process, allowed it to excel in our specific application, demonstrating the power and effectiveness of transfer learning in modern NLP tasks.

### 3.7.5 Results and Conclusion

By leveraging transfer learning with GPT-2, we achieved superior results in our project. The pre-trained model provided a strong foundation, and the fine-tuning process allowed us to adapt it to our specific needs, resulting in enhanced performance. This approach not only saved time and computational resources but also highlighted the effectiveness of transfer learning in tackling complex NLP tasks.

In conclusion, transfer learning, especially with advanced models like GPT-2, offers a powerful method for improving model performance with limited data. Our experience with GPT-2 underscores its potential and effectiveness in achieving state-of-the-art results in various applications.

# Chapter 4

## Discussion and Results

### 4.1 Effectiveness of LLMs

LLMs have revolutionized the field of NLP, enabling significant advancements in tasks such as machine translation, summarization, and title generation. The Transformer architecture, introduced by Vaswani et al. in the seminal paper "Attention is All You Need" [1], has been particularly influential. This architecture, based solely on attention mechanisms, has demonstrated superior performance compared to recurrent and convolutional neural networks [6, 7].

The core innovation of the Transformer is its use of self-attention mechanisms, which allow the model to weigh the importance of different words in a sequence relative to each other. This parallelizable architecture significantly reduces training time while improving the quality of the output. In our study, we leveraged this architecture for the task of generating titles for research papers based on their abstracts [8, 9].

### 4.2 LLM Hallucinations and Other Factors

One notable challenge encountered during the deployment of LLMs is the phenomenon known as "hallucination." Hallucination refers to the generation of plausible-sounding but factually incorrect or nonsensical outputs by the model. This issue arises due to the probabilistic nature of LLMs, where the model generates text based on learned patterns rather than factual accuracy [10, 3].

In our experiments, hallucinations manifested in the form of titles that, while grammatically correct and contextually relevant, did not accurately reflect the content



of the abstracts. Addressing this challenge requires incorporating strategies such as reinforcement learning from human feedback (RLHF) and refining the training data to emphasize factual correctness [7, 11].

Additionally, factors such as dataset quality, model architecture, and training procedures play critical roles in the performance of LLMs. Ensuring a high-quality, representative dataset is crucial for training models that generalize well to the intended task. Architectural innovations and advanced training techniques further enhance the capabilities of LLMs [4, 12].

### 4.3 Results

The progression of our experiments highlighted the transformative impact of model architecture and fine-tuning on performance. Initially, the custom-built LLM from scratch demonstrated limited effectiveness, with generated titles often lacking coherence and relevance.

### 4.4 Baseline LLM Performance

The performance of the baseline LLM, developed from scratch, was observed through qualitative analysis. The generated titles frequently exhibited issues such as incomplete phrases, irrelevant content, and lack of contextual understanding. For instance, when provided with the abstract of the paper "Attention Is All You Need," the custom-built LLM produced the title:

**Custom-Built LLM Title:** "Tranformer Model Sequence Tranduction Good Performance"

While this title captures some aspects of the paper, it lacks the specificity and context provided in the abstract.

### 4.5 GPT-2 Performance

Introducing GPT-2 into our experiments marked a significant improvement. The pre-trained GPT-2 model, when applied to the title generation task, produced more coherent and contextually appropriate titles. However, the performance was still

suboptimal in certain cases. For example, given the same abstract, GPT-2 generated the title:

**GPT-2 Title:** "The Transformer: A New Simple Network Architecture for Sequence Transduction"

This title is more detailed and relevant but still not fully reflective of the abstract's emphasis on attention mechanisms and the specific results of the study [8].

## 4.6 Fine-Tuned GPT-2 Performance

Fine-tuning GPT-2 on our specific dataset of research paper abstracts and titles resulted in a notable performance boost. The fine-tuned model generated titles that were not only coherent and grammatically correct but also accurately reflected the content and context of the abstracts. For the same abstract, the fine-tuned GPT-2 produced the title:

**Fine-Tuned GPT-2 Title:** "A Transformer Architecture for Sequence Transduction"

This title effectively captures the core contributions and findings of the paper, highlighting the use of attention mechanisms and the superiority of the Transformer architecture [9].

## 4.7 Case Study: Generated Titles

To illustrate the improvements in our title generation process, consider the following case studies. Each example highlights titles generated by different models based on the abstract of a research paper, as illustrated in Table 4.1.

The progression from the custom-built LLM to the fine-tuned GPT-2 demonstrates a significant enhancement in the quality and relevance of the generated titles [7, 11].

## 4.8 Conclusion

The effectiveness of LLMs in generating research paper titles is significantly influenced by the model architecture and the extent of task-specific fine-tuning. While initial models may struggle with coherence and relevance, leveraging pre-trained models like

GPT-2 and fine-tuning them on specific datasets can yield substantial improvements. Addressing challenges such as hallucinations and ensuring high-quality training data are critical for optimizing the performance of LLMs in specialized tasks [10, 6].

	<b>Custom-Built LLM Title</b>	<b>GPT-2 Title</b>	<b>Fine-Tuned GPT-2 Title</b>	<b>GPT-4</b>
Abstract #1	Tranformer Model Seqence Tranduction Good Performance	The Trans-former: A New Simple Network Architecture for Sequence Transduc-tion	A Trans-former Architecture for Sequence Transduc-tion	Attention Is All You Need
Abstract #2	Research Papper Title Transformer and Other Models with Good Result	Using Trans-former Mod-els to Gen-erate Titles for Research Papers: Initial Findings	Applying Transformer-Based Mod-els for Automated Research Paper Title Generation	Enhancing Research Paper Title Genera-tion with Fine-Tuned Transformer Models
Abstract #3	Machine Learnings Categorize Dark Web VPN TOR Traffic with High Acu-racy 3	Distinguishing VPN and TOR Traffic on the Dark Web Using Machine Learning	Applying Machine Learning Algorithms to Identify and Categorize VPN and TOR Traffic on the Dark Web	Utilizing Machine Learning for Accurate Classifica-tion of VPN and TOR Traffic on the Dark Web

Table 4.1: Comparison of Titles Generated by Different Models

# Chapter 5

## Conclusion

### 5.1 Summary of Findings

This project set out to explore the effectiveness of LLMs, particularly the Transformer architecture, in the task of generating titles for research papers based on their abstracts. Our investigation involved the development and evaluation of various model configurations, including a custom-built LLM, the pre-trained GPT-2 model, and a fine-tuned version of GPT-2.

The custom-built LLM from scratch demonstrated fundamental challenges, such as generating titles that often lacked coherence and relevance. The pre-trained GPT-2 model showed marked improvements due to its extensive training on diverse datasets, but it still struggled with fully understanding the specific context of research paper abstracts [8, 6]. Fine-tuning GPT-2 on our specific dataset of research paper abstracts and titles resulted in substantial performance enhancements, producing titles that were both coherent and contextually accurate.

### 5.2 Effectiveness of the Transformer Architecture

The Transformer architecture, with its self-attention mechanism, has proven to be highly effective for sequence transduction tasks. This architecture's ability to handle long-range dependencies and its parallelizable nature make it particularly suitable for tasks that require understanding complex relationships within the data. Our findings align with the broader literature, which has highlighted the superiority of Transformers over traditional recurrent and convolutional neural networks in many

NLP tasks [1, 7].

### 5.3 Addressing LLM Hallucinations and Challenges

One of the critical challenges in deploying LLMs is addressing the issue of hallucinations, where the model generates text that is plausible but factually incorrect or nonsensical. Our experiments highlighted this issue, particularly with the custom-built LLM and the pre-trained GPT-2 model [12, 10]. Fine-tuning and incorporating strategies such as reinforcement learning from human feedback (RLHF) can mitigate these issues, but they remain an area for further research [11, 2].

Additionally, the quality of the dataset, model architecture, and training procedures significantly influence the performance of LLMs. Ensuring high-quality, representative datasets and employing advanced training techniques are essential for optimizing model performance [4].

### 5.4 Implications for Future Research

Our research has several implications for future work in the field of NLP and beyond:

- **Fine-Tuning Pre-Trained Models:** Our results underscore the importance of fine-tuning pre-trained models on task-specific datasets. This approach can significantly enhance the performance of LLMs in specialized tasks [9].
- **Mitigating Hallucinations:** Addressing hallucinations remains a critical area for further research. Developing robust methods to ensure factual accuracy in generated text is crucial for the broader adoption of LLMs in applications requiring high reliability [12].
- **Expanding Application Domains:** The success of the Transformer architecture in various NLP tasks suggests potential applications in other domains, such as computer vision, speech recognition, and beyond. Exploring these possibilities can lead to new breakthroughs in artificial intelligence [7, 6].

## 5.5 Final Remarks

The transformative impact of LLMs, particularly the Transformer architecture, on the field of NLP is undeniable. Our study contributes to this growing body of knowledge by demonstrating the effectiveness of these models in generating research paper titles. By addressing the challenges and leveraging the strengths of LLMs, we can unlock new possibilities for intelligent text generation and beyond [1, 8].

In conclusion, the journey of exploring LLMs has been both challenging and rewarding. The insights gained from this research provide a solid foundation for future endeavors in the field, paving the way for continued advancements in artificial intelligence and machine learning.

# Bibliography

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in neural information processing systems*, 2017, pp. 5998–6008.
- [2] N. Patwardhan, S. Marrone, and C. Sansone, “Transformers in the real world: A survey on nlp applications,” *Information*, vol. 14, no. 4, p. 242, 2023.
- [3] A. Chernyavskiy, D. Ilvovsky, and P. Nakov, “Transformers: The end of history for nlp?” *arXiv preprint arXiv:2105.00813*, 2021.
- [4] A. Rahali and M. A. Akhloufi, “End-to-end transformer-based models in textual-based nlp,” *AI*, vol. 4, no. 1, pp. 54–110, 2023.
- [5] Falah.G.Salieh, “Research paper dataset 2023,,” 2023. [Online]. Available: Falah/research\_paper2023
- [6] X. Qiu, T. Sun, Y. Xu, Y. Shao, N. Dai, and X. Huang, “Generative pre-trained transformer: A comprehensive review on enabling technologies, potential applications, emerging challenges, and future directions,” *arXiv preprint arXiv:2005.10435*, 2020.
- [7] S. Islam, H. Elmekki, A. Elsebai, J. Bentahar, N. Drawel, and G. Rjoub, “A comprehensive survey on applications of transformers for deep learning tasks,” *arXiv preprint arXiv:2306.07303*, 2023.
- [8] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, “Language models are unsupervised multitask learners,” OpenAI, 2019.
- [9] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *arXiv preprint arXiv:2005.14165*, 2020.



- [10] T. Xiao and J. Zhu, “Introduction to transformers: an nlp perspective,” *arXiv preprint arXiv:2311.17633*, 2023.
- [11] Authors, “Tigen – title generator based on deep nlp transformer model for scholarly literature,” *SpringerLink*, 2023.
- [12] Y. Chen and S. Eger, “Transformers go for the lols: Generating (humorous) titles from scientific abstracts end-to-end,” in *Proceedings of the 4th Workshop on Evaluation and Comparison of NLP Systems*. Association for Computational Linguistics, 2023, pp. 62–84.