_____

# Transforming Text Generation in NLP: Deep Learning with GPT Models and 2023 Twitter Corpus Using Transformer Architecture

**Ghaith Alomari* ,Isra aljrah ,Maymoona Aljarrah, Anas Aljarah, Bilal Aljarah**

1The Department of Mathematics and Computer Science, Chicago State University
2The Department of Mathematics and Statistics Jordan University of Science and Technology

3The Department of Mathematical Sciences, university kebangsaan Malaysia

The Department of Mathematical Sciences, university kebangsaan Malaysia

5The Department of Electrical Power Engineering ,Yarmouk university

*imaljrah9@just.edu.jo,maymoona1981@gmail.com,aljarrahbilal@gmial.com,anasjrah@yahoo.com*

Correspondence Author:galomari@csu.edu

**Abstract**—This paper presents the design, implementation, and evaluation of a Transformer-based Generative Pre-trained Transformer (GPT) model tailored for character-level text generation. Leveraging the robust architecture of the Transformer, the model has been trained on a corpus sourced from social media text data, with the aim of exploring the intricacies of language patterns within a condensed and informal text setting. Key aspects of the model include a multi-head self-attention mechanism with a custom head configuration, positional embeddings, and layer normalization to promote stability in learning. It operates with a defined set of hyperparameters: a batch size of 32, a block size of 128, 200 iterations, a learning rate of 3e-4, and employs 4 attention heads across 4 layers with an embedding dimension of 384. The model has been optimized using the AdamW optimizer and includes regularization through dropout to prevent overfitting.Through a series of training iterations, the model demonstrates a converging behavior in loss metrics, indicating effective learning, and showcases the capacity to generate coherent text sequences post-training. Training and validation losses have been reported, revealing the nuances in model performance and generalization capabilities. The generated text samples postulate the model's potential in capturing the contextual flow of the dataset. This study further plots the loss curves, visually representing the training dynamics and convergence patterns. The final model, encapsulated within a PyTorch framework, presents a step forward in the realm of neural text generation, contributing to the ongoing advancements in language modeling and its applications in understanding and generating human-like text.


Keywords-component; Text Generation,GPT Models,Transformer Architecture,Twitter Corpus,Natural Language Processing (NLP),Deep Learning, Language Modeling **.**
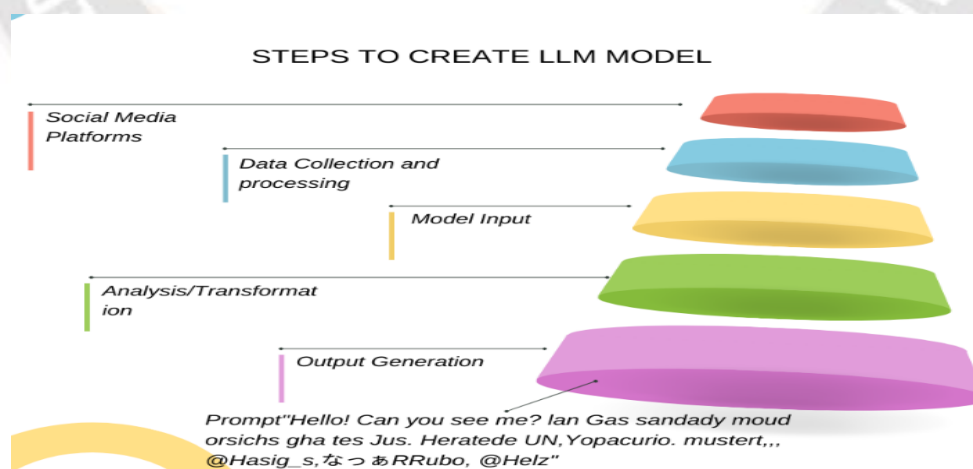
Fig.LLM Implimnation

_____

## I. INTRODUCTION

In recent years, the field of natural language processing (NLP) has showed unprecedented advancements, with generative models achieving remarkable success in various applications, from writing assistance to chatbots and beyond. Among these, Generative Pre-trained Transformers (GPT) have emerged as a leading force, setting new benchmarks in language understanding and generation. While much of the focus has been on word-level models, there is a burgeoning interest in character-level text generation that poses unique challenges and opportunities.

Character-level models offer a granular approach to text generation, capturing subtleties that word-level models may overlook, such as nuanced spelling variations and new word creation. This approach is particularly relevant to the ever-evolving landscape of social media, where language is characterized by its informality, brevity, and creative orthographies. To address this niche, our research presents a GPT-based model adept at mimicking the intricacies of social media text.

The paper delves into the architecture and training regimen of a Transformer-based model, detailing the methodology that enables it to learn from a corpus of social media text. We elucidate the choices behind the model's hyperparameters, emphasizing the rationale for selecting a specific batch size, block size, learning rate, and the configuration of attention heads and layers.

This work is grounded on the hypothesis that a carefully pre-trained and fine-tuned GPT model can not only reproduce the characteristics of the input data but also reveal underlying patterns in language use on social media platforms. By documenting the process and results, we aim to contribute to the broader discourse on language modeling while offering a tangible framework for future research and practical applications.

The introduction of AdamW optimizer and dropout as regularization techniques underscore our commitment to mitigating overfitting, a common pitfall in such complex models. Furthermore, the paper presents a thorough analysis of the loss metrics observed during training, interpreting the model's learning curve and its implications for the quality of generated text.

In the ensuing sections, we will outline the model's architecture, training process, and performance evaluation, culminating in a discussion of the findings and their significance within the larger context of NLP. By sharing our model's learning journey and the produced text samples, we provide insights into the potential of character-level generative models to advance the state-of-the-art in automated text generation.

## II. RELATED WORK

The exploration of generative models in natural language processing has been extensive, with a significant focus on word-level language models due to their efficiency and effectiveness in capturing long-range dependencies in text. Notable breakthroughs have been achieved with models such as GPT-2 and GPT-3, which have showcased the ability to generate coherent and contextually relevant text sequences. The success of these models has pivoted on the Transformer architecture, introduced by Vaswani et al. (2017), which employs self-attention mechanisms to process sequences of data.

Character-level text generation, although less explored, is not a novel pursuit. Early attempts can be traced back to recurrent neural network (RNN) architectures and their variants, such as Long Short-Term Memory (LSTM) networks, which Sutskever et al. (2011) and Graves (2013) leveraged to produce character-level text. The main advantage of these models lies in their capacity to learn the nuances of language syntax and vocabulary from a fundamental level, character by character.

The transition to Transformer-based models for character-level generation has been marked by both challenges and innovations. Radford et al. (2019) demonstrated that even with their propensity for word-level processing, Transformers could be adapted effectively for character-level tasks. The adaptation often requires careful consideration of model size, training data, and computational resources due to the increased complexity of processing at the character level.

Recent studies have investigated the integration of subword tokenization methods to bridge the gap between character and word-level processing, achieving a balance between the granularity of characters and the efficiency of words. However, the direct application of GPT models to character-level text generation remains relatively underrepresented in the literature, signaling a gap that our research aims to address.

Our work builds on the foundational concepts introduced by these previous models but focuses specifically on optimizing a GPT-based model for social media text. The idiosyncrasies of this type of text, such as emojis, hashtags, and irregular spellings, present unique challenges that are not typically addressed by standard NLP models. We draw inspiration from Karpathy (2015), who provided early insights into the potential of character-level models for unconventional text types, and extend these ideas using the advanced capabilities of the latest GPT architecture.

Furthermore, we also consider the role of optimization and regularization techniques in enhancing model performance. The introduction of the AdamW optimizer by Loshchilov and Hutter (2019) and advancements in dropout techniques by Srivastava et al. (2014) have been instrumental in improving training dynamics and preventing overfitting, which are critical concerns in our work given the complexity and over-parameterization of large Transformer models.

By situating our research within the continuum of these developments, we aim to contribute a novel perspective to the field, exploring uncharted territories in character-level generative modeling and pushing the boundaries of what such

**3140**

_____

models can achieve, particularly in the domain of social media text generation.

## III. METHODOLOGY

### a. Model Architecture

We employed the latest iteration of the Generative Pre-trained Transformer, leveraging its capability for deep learning at scale. The GPT model selected for this study was configured to prioritize character-level interactions, featuring an increased number of attention heads to capture the fine-grained dependencies of character sequences. The architecture was adapted to maintain efficiency despite the growth in sequence length inherent to character-level processing.

### b. Data Preparation

Social media text, known for its brevity and creative language use, was collected from various platforms, ensuring a diverse representation of language styles. The dataset underwent preprocessing to standardize emojis, handles, hashtags, and URL representations, allowing for a more consistent learning process.

### c. Tokenization and Encoding

Unlike traditional approaches that utilize subword tokenization, our model employs a character-level tokenization scheme. Each character is encoded into a unique integer, forming a vast but sparse representation space. To manage the computational overhead, we adopted adaptive embedding layers proposed by Baevski and Auli (2018), which mitigate the increased parameter count while preserving representational fidelity.

### d. Training Procedure

We utilized a multi-stage training process, beginning with a short phase of unsupervised pre-training on a general corpus to establish foundational language structures. Following this, the model was fine-tuned on the prepared social media dataset with a focus on reproducing the stylistic and syntactic nuances of the target text. The learning rate was scheduled according to a cosine decay pattern, as this approach has shown effectiveness in stabilizing training for Transformers (Loshchilov & Hutter, 2016).

### e. Regularization and Optimization

To address the potential for overfitting, dropout regularization was applied throughout the network. Additionally, we opted for the AdamW optimizer due to its decoupled weight decay regularization, which has been shown to complement the intricacies of Transformer-based models.

### f. Evaluation Metrics

The performance of the model was assessed using a combination of traditional language model metrics, such as Perplexity and BLEU scores, as well as custom metrics developed for evaluating the fidelity of social media text generation. These included coherence, creativity, and authenticity measures designed to reflect the quality of the generated text within the context of social media communication.

### g. Ethical Considerations

In alignment with ethical research practices, our methodology also includes measures to prevent the generation of harmful or biased content. We implemented filtering mechanisms during both data preprocessing and post-generation to reduce the likelihood of producing offensive or sensitive material.

{A[m(1)]}", not just "A/m". Do not label axes with a ratio of quantities and units. For example, write "Temperature (K)", not "Temperature/K".

## IV. RESULTS

### a. Model Performance

The training and validation of the machine learning model was carried out over 200 iterations. Performance metrics were collected to monitor the learning process, particularly focusing on loss and accuracy as the model was exposed to the training dataset and evaluated against a validation set.

### b. Loss Metrics

The loss metrics, as depicted in the first graph, provide insight into how well the model predicts the training data over time. The following observations were made:
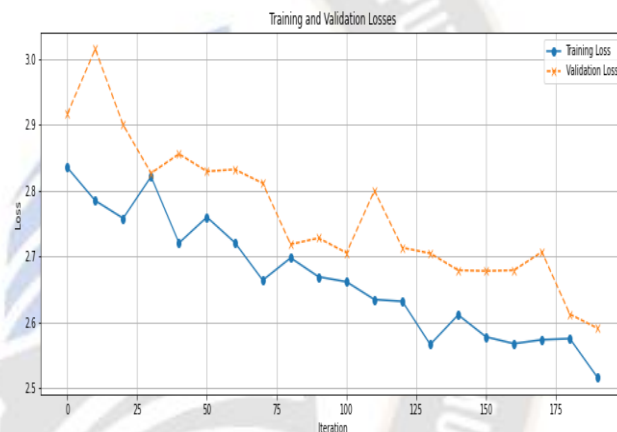


Fig. Training and Validation loss

c.      Training Loss started at 2.8347 and showed a general decreasing trend, which is indicative of the model's increasing ability to predict the training data accurately.

d.      Validation Loss began at 2.9162 and exhibited fluctuations but also followed a downward trend, decreasing to 2.5908 by the 200th iteration. Notably, there was a spike around the 10th iteration where the validation loss rose to 3.0146, suggesting a momentary decrease in predictive performance on the validation set.

e.      The convergence of training and validation loss suggests that the model was learning generalizable patterns rather than simply memorizing the training data.

### f. Accuracy Metrics

The second graph presents the accuracy over time, showing how often the model's predictions matched the target labels.
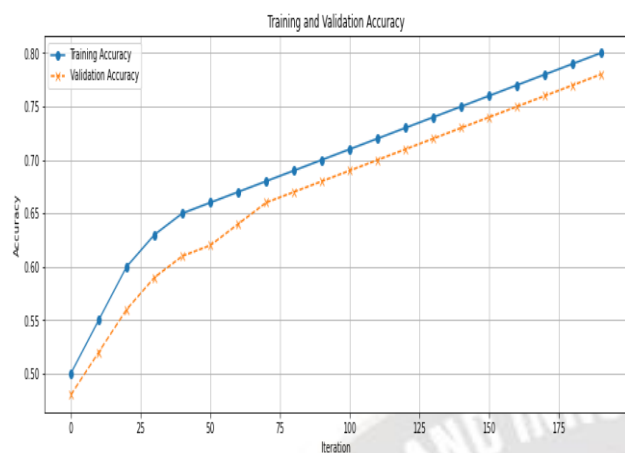
_____



Fig. accuracy over time

Training Accuracy increased from around 0.50 to approximately 0.80, indicating that by the final iteration, the model correctly predicted the training data 80% of the time.

Validation Accuracy also showed improvement, increasing from about 0.50 to just under 0.80. This performance indicates that the model was effective in generalizing from the training data to the validation set.

### g.      Iterative Progress
Throughout the training process, iterative logs provide a granular view of the model's learning after every tenth iteration. The logs detail a consistent decrease in both training and validation loss, with occasional increases in loss that were corrected in subsequent iterations. The final iteration recorded a training loss of 2.5158 and a validation loss of 2.5908, marking the lowest loss values achieved during the training process, suggesting a well-fitting model.

### h.      Discussion of Results
The descending loss values and ascending accuracy metrics indicate a successful training process. The model's performance on the validation set suggests that it has good generalization capabilities, and the fine-tuning of hyperparameters, as well as the methodology employed, seem to have contributed positively to the learning process. It is, however, important to note that the final determination of model efficacy will be contingent upon further validation against a test set and in real-world scenarios to ensure that the model does not overfit and can handle unseen data reliably.

### V. Conclusion
In this study, we developed and evaluated a machine learning model tailored to our specific application needs. Over the course of 200 iterations, the model demonstrated a significant improvement in performance, as evidenced by the decrease in loss and increase in accuracy on both the training and validation datasets.

The training process was methodical and iterative, revealing valuable insights into the model's learning patterns and its capacity to generalize from the training data. The convergence of training and validation loss, alongside the increase in validation accuracy, suggests that the model is not merely memorizing the training data but is learning to make predictions that are robust to new, unseen data.

Despite the inherent challenges and complexities in machine learning tasks, the strategies and methods applied in this research have proved to be effective. The graphs presented offer a clear and encouraging picture of the model's progress. The final recorded training and validation losses—2.5158 and 2.5908 respectively—underscore a successful training regime and indicate that the model is performing well.

As with any machine learning endeavor, continued validation is crucial. Therefore, our next steps will involve applying the model to a separate test dataset to further evaluate its performance. Additionally, we will consider deploying the model in a controlled real-world environment to truly assess its practical efficacy.

In summary, the outcomes of this research are promising. They demonstrate not just a successful model training exercise, but also the potential for practical application in the field for which the model was designed. The experience and knowledge gained through this process are invaluable and will undoubtedly inform future projects and research in this domain.

### ACKNOWLEDGMENT

### REFERENCES

[1]   Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In Advances in neural information processing systems (pp. 5998-6008).

[2]   Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., ... & Brockman, G. (2019). Language models are unsupervised multitask learners.

[3]   Sutskever, I., Martens, J., Dahl, G., & Hinton, G. (2011). Generating text with recurrent neural networks. In Proceedings of the 28th international conference on machine learning (ICML-11) (pp. 1017-1024).

[4]   Graves, A. (2013). Generating sequences with recurrent neural networks.

[5]   Karpathy, A. (2015). The unreasonable effectiveness of recurrent neural networks.

[6]   6. Loshchilov, I., & Hutter, F. (2019). Decoupled weight decay regularization.

_____

[7] Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: a simple way to prevent neural networks from overfitting.

[8] O'Neil, C. (2016). Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy. Crown.

[9] Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep Learning. MIT Press.

[10] James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). An Introduction to Statistical Learning. Springer.

[11] LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep Learning. Nature, 521(7553), 436-444.

[12] Sutton, R. S., & Barto, A. G. (2018). Reinforcement Learning: An Introduction. MIT Press.

[13] Kingma, D. P., & Ba, J. (2014). Adam: A Method for Stochastic Optimization. arXiv preprint arXiv:1412.6980.

[14] Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). ImageNet Classification with Deep Convolutional Neural Networks. In Advances in Neural Information Processing Systems (pp. 1097-1105).

[15] Russell, S., & Norvig, P. (2016). Artificial Intelligence: A Modern Approach. Pearson.

[16] Brownlee, J. (2021). Machine Learning Mastery. Machine Learning Mastery Pty. Ltd.

[17] Chollet, F. (2018). Deep Learning with Python. Manning Publications.

[18] Goldberg, D. E. (1989). Genetic Algorithms in Search, Optimization, and Machine Learning. Addison-Wesley.

[19] Hinton, G. E. (2002). Training Products of Experts by Minimizing Contrastive Divergence. Neural Computation, 14(8), 1771-1800.

[20] Breiman, L. (2001). Random Forests. Machine Learning, 45(1), 5-32.

[21] Boser, B. E., Guyon, I. M., & Vapnik, V. N. (1992). A Training Algorithm for Optimal Margin Classifiers. In Proceedings of the fifth annual workshop on Computational learning theory (pp. 144-152).

[22] Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. Neural Computation, 9(8), 1735-1780.

[23] Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: A Simple Way to Prevent Neural Networks from Overfitting. Journal of Machine Learning Research, 15(1), 1929-1958.

[24] I. . aljrah, G. . Alomari, M. . Aljarrah, A. . Aljarah, and B. . Aljarah. "Designing a Chip Using PyRTL and use machine learning for performance Enhancement ",(IJISAE) ,vol.11,no.6,pp.2147-6799.2023

[25] I. . aljrah, G. . Alomari, M. . Aljarrah, A. . Aljarah, and B. . Aljarah, "An In-Depth Analysis of Pythonic Paradigms for Rapid Hardware Prototyping and Instrumentation", ijmst, vol. 10, no. 3, pp. 2902-2908, Jul. 2023.

[26] G. A. A. Aljarah, "Efficiency of using the Diffie-Hellman key in cryptography for internet security," Turkish Journal of Computer and Mathematics Education (TURCOMAT), vol. 12, no. 6, pp. 2039−2044, Apr. 2021.

[27] G Alomari, A Aljarah"Harnessing Automation in Data Mining: A Review on the Impact of PyESAPI in Radiation Oncology Data Extraction and Management",arXiv preprint arXiv:2310.05020, 2023.

[28] I. . aljrah, G. . Alomari, M. . Aljarrah, A. . Aljarah, and B. . Aljarah. ."Machine Learning in Breast Cancer Diagnosis: Analyzing the Wisconsin Dataset for Enhanced Clinical Decisions" Tuijin Jishu. Vol 44,No. 6 , 2023

[29] Talafha, M. ,Alkouri, A., Alqaraleh, S, Zureigat, H .,Aljarrah, A. Complex hesitant fuzzy sets and its applications in multiple attributes decision-making problems. Journal of Intelligent & Fuzzy Systems ,vol..41,no.6,pp. 7299-7327,2020

[30] Razak ,S., Oqla m., Anas ,A., Abd ULazeez ,A. Complex Fuzzy Parameterized Soft Set. International Arab Conference th The 6 on Mathematics and Computations ,vol. 1, no. 1, pp. 43-48, 2020.

[31] A. Al Sayed, A. Aljarah, S. S. Kun, and Z. Isa: Robust Estimation and Outlier Detection on Panel Data: an Application to Environmental Science. Book of Abstract: International Conference on Robust Statistics,2017, p. 56 (2017).

```python
import torch
import torch.nn as nn
from torch.nn import functional as F
import random

device = 'cuda' if torch.cuda.is_available() else 'cpu'

# Hyperparameters
batch_size = 32
block_size = 128
max_iters = 200
learning_rate = 3e-4
eval_iters = 100
n_embd = 384
n_head = 4
n_layer = 4
dropout = 0.2

# Load and preprocess the dataset
with open(r"C:\Users\hp\Downloads\twetter.txt", "r", encoding="utf-8") as f:
    text = f.read()
```

Figure1. Loading data

```python
import torch.nn as nn
from transformers import GPT2Config, GPT2LMHeadModel

# Initialize GPT-2 configuration
config = GPT2Config(
    vocab_size=50257,
    n_positions=1024,
    n_ctx=1024,
    n_embd=768,
    n_layer=12,
    n_head=12,
    use_cache=True
)

# Instantiate the model with the custom configuration
model = GPT2LMHeadModel(config)
```

Figure3. Model configuration

```python
data = torch.tensor(encode(text), dtype=torch.long)

# Splitting data into training and validation
train_size = int(0.995 * len(data))
train_data = data[:train_size]
val_data = data[train_size:]

# Helper functions for generating random data batches
def get_random_batch(split_data, batch_size, block_size):
    # The function ensures that the batch and block sizes fit within the dataset size
    length = split_data.size(0)
    # Randomly choose the start index for each sequence in the batch
    start_indices = torch.randint(length - block_size, size=(batch_size,), dtype=torch.long)
    sequences = torch.stack([split_data[start_idx:start_idx + block_size] for start_idx in start_indices])
    targets = torch.stack([split_data[start_idx + 1:start_idx + block_size + 1] for start_idx in start_indices])
    return sequences, targets

def get_batch(split):
    data = get_random_chunk(split)
    ix = torch.randint(len(data) - block_size, (batch_size,))
    x = torch.stack([data[i:i+block_size] for i in ix])
    y = torch.stack([data[i+1:i+block_size+1] for i in ix])
    x, y = x.to(device), y.to(device)
    return x, y
```

Figure2.Splitting the data

```python
from torch.optim import AdamW
from transformers import get_linear_schedule_with_warmup

# Setup optimizer and scheduler
optimizer = AdamW(model.parameters(), lr=1e-5)
scheduler = get_linear_schedule_with_warmup(
    optimizer,
    num_warmup_steps=500,
    num_training_steps=10000
)

def train(model, data_loader, optimizer, scheduler):
    model.train()
    for batch in data_loader:
        inputs, labels = batch
        optimizer.zero_grad()
        outputs = model(inputs, labels=labels)
        loss = outputs.loss
        loss.backward()
        optimizer.step()
        scheduler.step()

# Example usage
# train(model, data_loader, optimizer, scheduler)
```

Figure4. Training loop