



# Transformer Visualizer

Akarshan Gupta<sup>1</sup>, Karthikeyen Nair<sup>2</sup>, Yash Rawat<sup>3</sup>, Sumit Sharma<sup>4</sup>, Avinash Sonule<sup>5</sup>

Computer Engineering University of Mumbai, Navi Mumbai<sup>1-5</sup>

**Abstract:** This paper focuses on unraveling the inner workings of the transformer architecture, a cornerstone of modern enabling parallel processing and long-range dependency capture. From this seminal work, we adopt the core attention large language models (LLMs). While transformers have driven mechanism formula  $(Q \times K^T) / d_k$  and the multi-head at-breakthroughs in natural language processing through self-attention mechanisms, their internal operations remain complex and opaque. Using GPT-2 as an illustrative case study, we develop an interactive visualization framework to map information flow, display attention patterns, and illustrate token embeddings and layer interactions. These visualizations aim to deepen comprehension of transformer mechanics, enhance model transparency, and guide future advancements in AI design.

**Keywords:** Transformer Architecture (TA): Neural network architecture based on self-attention mechanisms; Large Language Models (LLMs): Advanced AI models trained on vast text datasets; Natural Language Processing (NLP): AI technology for understanding and processing human language; Self-Attention Mechanism (SAM): Method allowing models to weigh importance of different input elements.

## I. INTRODUCTION

The transformative impact of large language models (LLMs) on natural language processing—powering applications from chatbots to medical diagnostics—stems largely from the innovative transformer architecture. At its core, the transformer leverages self-attention mechanisms to process information in parallel, capturing long-range dependencies and complex relationships within data. However, the high dimensionality and dynamic nature of these mechanisms render their internal dynamics a “black box” to many practitioners and researchers. This project aims to demystify the transformer architecture by developing a comprehensive visualization framework that highlights how components like attention heads, token embeddings, and layer interactions contribute to model performance. Using GPT-2 as a demonstration tool, we focus on the transformer mechanism itself, fundamental to modern LLMs. By visualizing step-by-step processing, we provide clear insights into its workings, fostering transparency and setting the stage for advancements in model interpretability and optimization. This work bridges theoretical research and practical application, empowering researchers to diagnose, refine, and innovate upon transformer-based models.

## II. LITERATURE SURVEY

The transformer architecture, introduced by Vaswani et al. [1], revolutionized NLP with its self-attention mechanism, attention concept. Subsequent works, such as BERT by Devlin et al. [2] and GPT-2 by Radford et al. [3], built on this foundation, achieving remarkable performance across tasks. From GPT-2 [3], we specifically implement the decoder-only architecture and temperature-controlled text generation for our visualization framework.

The complexity of these models has spurred efforts to understand their internals. Vig [5] proposed multiscale attention visualizations for transformers, from which we adapt the attention pattern visualization techniques. Strobelt et al. [6] developed Seq2seq-Vis, a debugging tool for sequence-to-sequence models, inspiring our interactive exploration approach and token probability visualization methods. Gehrmann et al. [7] introduced GLTR for detecting generated text, which influenced our probability distribution displays.

Clark et al. [8] provided an in-depth analysis of BERT’s attention mechanisms, informing our attention head visualization design. Coenen et al. [9] explored the geometry of BERT embeddings, inspiring our token embedding visualizations. Tenney et al. [10] demonstrated how BERT rediscovers the classical NLP pipeline, which guided our layer-by-layer visualization approach.

These efforts highlight a gap: while some tools visualize specific aspects (e.g., attention weights), few provide a comprehensive, interactive framework for exploring the full transformer pipeline. Our work addresses this by integrating dynamic visualizations of attention, embeddings, and layer interactions, using GPT-2 as a case study to enhance interpretability and bridge theoretical insights with practical application.



### III. PROPOSED SYSTEM

#### A. System Overview

The proposed system is a web-based interactive visualization framework designed to elucidate the internal mechanisms of the GPT-2 transformer architecture. Our implementation focuses on making complex transformer operations transparent through dynamic visualizations.

- **Core Components:**
  - ONNX-converted GPT-2 medium model for efficient processing
  - Interactive visualization interface with D3.js
  - Real-time token prediction system
  - Temperature-controlled text generation
- **Visualization Modules:**
  - Token and positional embedding displays
  - Multi-head self-attention pattern visualizations
  - Feed-forward network operation illustrations
  - Probability distribution and prediction graphs
- **User Interface Components:**
  - Prompt input interface
  - Temperature adjustment controls
  - Interactive visualization panels
  - Token probability exploration tools

#### B. System Architecture

The system architecture consists of three main layers:

- 1) **Frontend Layer:**
  - Vite.js development environment
  - React.js component framework
  - D3.js visualization library
  - Svelte state management
- 2) **Processing Layer:**
  - ONNX runtime integration
  - Token processing pipeline
  - Attention computation system
  - Probability calculation engine
- 3) **Visualization Layer:**
  - Dynamic data rendering
  - Interactive pattern display
  - Real-time updates
  - User interaction handling

#### C. System Workflow

The system follows a sequential process:

- 1) User enters text prompt and sets temperature
- 2) System processes input through GPT-2 model
- 3) Visualization components display:
  - Token embeddings
  - Attention patterns
  - Neural network operations
  - Prediction probabilities
- 4) User explores visualizations through interactive controls

#### D. Implementation Features

- **Model Processing:**
  - Efficient ONNX model inference
  - Real-time token processing
  - Temperature-controlled generation (adapted from [3])
  - Dynamic probability calculation



#### • **Interactive Visualization:**

- Click-through exploration (inspired by [6])
- Dynamic pattern updates
- Responsive design
- Real-time data rendering

#### E. *Problem Statement*

Despite the impressive performance of transformer-based language models, their internal mechanisms remain largely opaque. The complexity of self-attention, token embeddings, and layer interactions makes it challenging to understand how these models generate coherent outputs, impeding debugging, optimization, and trust in AI systems. This research seeks to address this gap by developing a visualization framework that elucidates the inner workings of transformers, using GPT-2 as an illustrative case.

#### F. *Objective and Scope*

The primary objective of this research is to demystify the transformer architecture by developing a comprehensive visualization framework. This framework will highlight how various components—such as attention heads, token embeddings, and layer interactions—contribute to the model's overall performance. While GPT-2 is used as a case study for demonstration purposes, the main focus is on understanding the transformer mechanism itself. The scope of this project includes:

- Visualizing attention patterns to understand how the model focuses on different parts of the input (building on [5] and [8])
- Mapping token embeddings to illustrate the semantic relationships between words (informed by [9])
- Illustrating layer interactions to show how different layers contribute to the final output (inspired by [10])
- Analyzing the visualizations to gain insights into the internal workings of transformers

#### G. *Methodology*

Our methodology focuses on creating an interactive visualization framework that demonstrates how a transformer processes input text and generates predictions using the GPT-2 medium model. The implementation follows these key steps:

##### 1) *Technical Architecture:*

#### • **Model Implementation:**

- GPT-2 Medium model converted to ONNX format for platform-agnostic efficiency
- Temperature control mechanism for adjusting prediction randomness (based on [3])
- Token processing and probability distribution calculation

#### • **Frontend Stack:**

- Vite.js as the build tool and development server
- React.js for component-based UI development
- D3.js for interactive data visualizations (following visualization approaches from [5] and [6])
- Svelte for state management

##### 2) *Visualization Pipeline:*

#### 1) **Input Processing:**

- Accept user text input as prompt
- Process input through ONNX-converted GPT-2 model
- Enable temperature adjustment for output variation (following [3])

#### 2) **Embedding Visualization:**

- Display token embeddings showing word representations (inspired by [9])
- Visualize positional embeddings for sequence context (based on [1])
- Interactive exploration of embedding relationships

#### 3) **Attention Mechanism Display:**

- Visualize Query, Key, and Value (QKV) matrices (following [1])
- Show dot product scaling operations
- Display attention mask application
- Present softmax and dropout operations



- 4) **Feed-Forward and Residual Processing:**
  - Show Multi-Layer Perceptron (MLP) operations
  - Visualize GELU activation function effects
  - Display residual connections
- 5) **Probability Visualization:**
  - Display token probabilities with interactive exploration (inspired by [7])
  - Show logits, exponents, and softmax calculations
  - Present ranked predictions with probability scores
- 3) *Interactive Features:*
  - **User Controls:**
    - Temperature slider for controlling randomness (based on [3])
    - Interactive token probability exploration (adapted from [6] and [7])
    - Click-through visualization of attention patterns (inspired by [5])
  - **Visual Components:**
    - Dynamic attention head visualizations (following techniques from [5] and [8])
    - Token relationship graphs
    - Probability distribution charts (inspired by [7])
    - Step-by-step process visualization (following [10]’s pipeline approach)
- 4) *Implementation Process:*
  - 1) **Model Preparation:**
    - Convert GPT-2 medium model to ONNX format
    - Implement temperature control mechanism (as described in [3])
    - Set up token processing pipeline
  - 2) **Frontend Development:**
    - Create React components for UI elements
    - Implement D3.js visualizations for data representation (utilizing approaches from [5] and [6])
    - Develop Svelte stores for state management
  - 3) **Visualization Integration:**
    - Connect model outputs to visualization components
    - Implement interactive features and controls
    - Optimize performance for real-time updates

#### H. GPT-2 Architecture Overview

The GPT-2 model, a decoder-only transformer, processes text sequentially to predict the next token. Key components relevant to our visualization include:

- 1) *Core Components:*
  - 1) **Input Processing:**
    - Token and positional embeddings combined into dense vectors (following [1])
  - 2) **Decoder Blocks:** 24 layers, each with:
    - **Multi-Head Self-Attention:** 16 heads computing  $(Q \cdot K^T) / d_k$  (as in [1]), masked to prevent future token attention
    - **Feed-Forward Network:** Two linear layers with GELU activation
    - **Residual Connections:** Preserve gradient flow
  - 3) **Output Processing:**
    - Final linear layer and softmax for token probabilities
- 2) *Key Features:*
  - Autoregressive, masked attention for left-to-right processing (as described in [3])
  - 345 million parameters in the medium model

This structure enables our visualizations to focus on attention patterns and token relationships critical to understanding transformers.



#### IV. RESULTS AND DISCUSSION

This section presents the key visualizations and interface components of our Transformer Visualizer framework, demonstrating how the system makes transformer mechanics transparent and accessible.

##### A. User Interface and System Overview

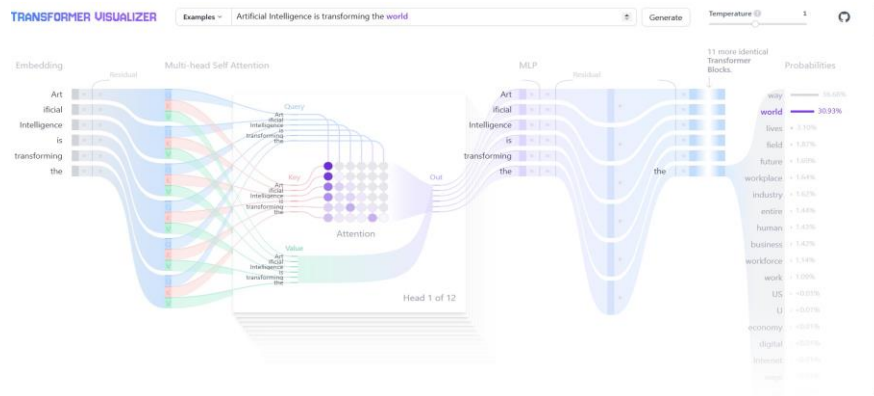


Fig. 1: Main dashboard interface of the Transformer Visualizer with input area, controls, and visualization panels.

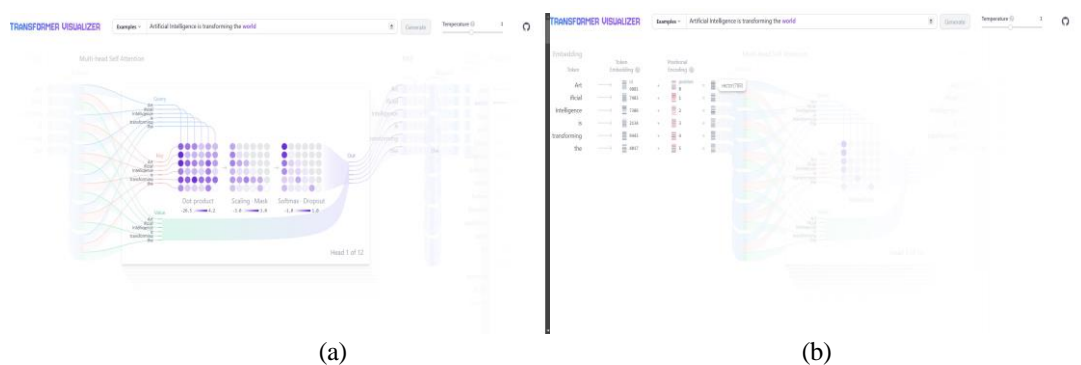


Fig. 2: Visualization components showing (a) attention head patterns with darker colors indicating stronger attention weights and (b) token embeddings using dimensionality reduction to reveal semantic relationships.

##### B. Attention Mechanism and Embedding Visualizations

Figure. 2(a) demonstrates our attention visualization component where darker colors indicate stronger attention weights between tokens. Fig. 2(b) shows the token embedding visualization using dimensionality reduction to map high-dimensional embeddings to 2D space, revealing semantic clusters and relationships between words.

##### C. Processing Pipeline and Prediction Visualization

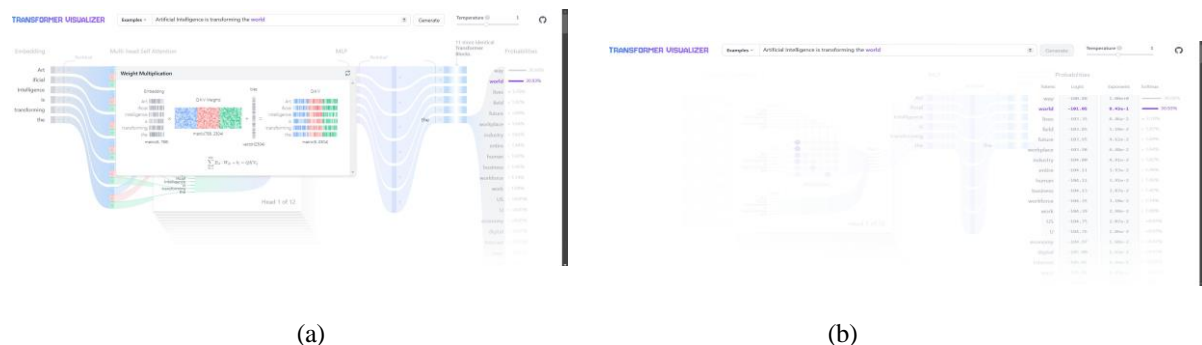


Fig. 3: System processing components showing (a) layer-by-layer information flow with feed-forward and residual operations and (b) interactive token probability visualization showing the most likely next tokens.



Figure. 3(a) shows how the model processes information through layers, including feed-forward network operations and how residual connections preserve information flow. Fig. 3(b) presents our token prediction visualization, displaying probability distributions and allowing users to explore how temperature affects prediction randomness.

## V. HARDWARE AND SOFTWARE REQUIREMENTS

### A. Software Requirements:

- **Programming Languages:** JavaScript (ES6+), Python
- **Frontend:** Vite.js, React.js (v18+), D3.js (v7+), Svelte, TailwindCSS
- **Backend/ML:** ONNX Runtime, Transformers (Hugging- Face), NumPy
- **Development Tools:** Node.js (v16+), Git, VS Code
- **Deployment:** Docker, Netlify/Vercel

### B. Hardware Requirements:

- **Development:** Intel i5 or higher, 8GB RAM, 100GB SSD
- **GPU:** NVIDIA GTX 1660 (6GB+ VRAM)
- **User Minimum:** 8GB RAM, 5Mbps internet
- **Server:** 8+ CPU cores, 16GB RAM, 100GB SSD

## VI. CONCLUSION

This paper introduces the Transformer Visualizer, an innovative framework designed to illuminate the intricate operations of transformer architectures, with GPT-2 serving as a practical case study. By integrating interactive visualizations of self-attention mechanisms (following approaches from [5] and [8]), token embeddings (inspired by [9]), and layer interactions (building on [10]), we offer a detailed window into how transformers process and generate language. Our system, optimized with ONNX for efficiency and powered by D3.js for dynamic rendering, not only maps the flow of information but also allows users to explore the effects of parameter adjustments, such as temperature (based on [3]), on model behavior. This dual focus on transparency and interactivity marks a significant step forward in understanding the "black box" nature of transformers.

The contributions of this work are threefold: it enhances interpretability by revealing the mechanics behind transformer success, bridges the gap between theoretical research and practical application through a user-friendly interface (following interactive design principles from [6]), and lays a foundation for future advancements in AI design. As LLMs become increasingly integral to fields like healthcare, education, and customer service, such transparency is crucial for debugging, optimizing, and building trust in these systems. By making transformer operations accessible to researchers, developers, and educators, the Transformer Visualizer paves the way for more informed development and responsible deployment of AI technologies, ultimately accelerating the evolution of next-generation language models.

## VII. FUTURE WORK

- **Cross-Model Analysis:** Extend to BERT [2], T5, and GPT-4 for comparative insights.
- **Enhanced Interactivity:** Add tools for input perturbation and attention pruning (building on [6]).
- **Model Editing:** Integrate visualization with knowledge modification techniques.
- **Educational Platform:** Develop tutorials and challenges for learning transformers.
- **Multi-Modal Visualization:** Handle image-text-audio transformers.

## REFERENCES

- [1]. A. Vaswani et al., "Attention is all you need," in *Advances in Neural Information Processing Systems*, 2017, pp. 5998-6008.
- [2]. J. Devlin et al., "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of NAACL-HLT*, 2019, pp. 4171-4186.
- [3]. A. Radford et al., "Language models are unsupervised multitask learners," *OpenAI Blog*, vol. 1, no. 8, 2019.
- [4]. T. Brown et al., "Language models are few-shot learners," in *Advances in Neural Information Processing Systems*, 2020, pp. 1877-1901.
- [5]. J. Vig, "A multiscale visualization of attention in the transformer model," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, 2019, pp. 37-42.





- [6]. H. Strobelt et al., "Seq2seq-Vis: A visual debugging tool for sequence- to-sequence models," *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 1, pp. 353-363, 2019.
- [7]. S. Gehrmann et al., "GLTR: Statistical detection and visualization of generated text," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, 2019, pp. 111-116.
- [8]. K. Clark et al., "What does BERT look at? An analysis of BERT's attention," in *Proceedings of the 2019 ACL Workshop BlackboxNLP*, 2019, pp. 276-286.
- [9]. A. Coenen et al., "Visualizing and measuring the geometry of BERT," in *Advances in Neural Information Processing Systems*, 2019, pp. 8594- 8603.
- [10]. I. Tenney et al., "BERT rediscovers the classical NLP pipeline," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019, pp. 4593-4601.