# Quantum Transformers with Variational Quantum Circuits for Enhanced Self-Attention in Machine Learning and Physics Simulations

**Arul Tripathi**

**Delhi, India**

arultripathi@gmail.com

## Abstract

The integration of quantum computing with classical machine learning architectures has opened new avenues for enhancing computational efficiency and model performance. In this work, we propose *Quantum Transformers*, a novel framework that leverages *Variational Quantum Circuits (VQCs)* to augment the self-attention mechanism in transformer models. By encoding classical data into quantum states and utilizing parameterized quantum gates, our approach enables the exploration of high-dimensional Hilbert spaces, offering a quantum-enhanced representation of attention weights. This hybrid quantum-classical architecture is designed to improve the expressivity and efficiency of self-attention, particularly in complex tasks such as natural language processing, time-series analysis, and physics simulations. Our results suggest that the synergy between quantum computing and transformer models holds significant promise for advancing both machine learning and computational physics, paving the way for scalable quantum-enhanced learning systems.
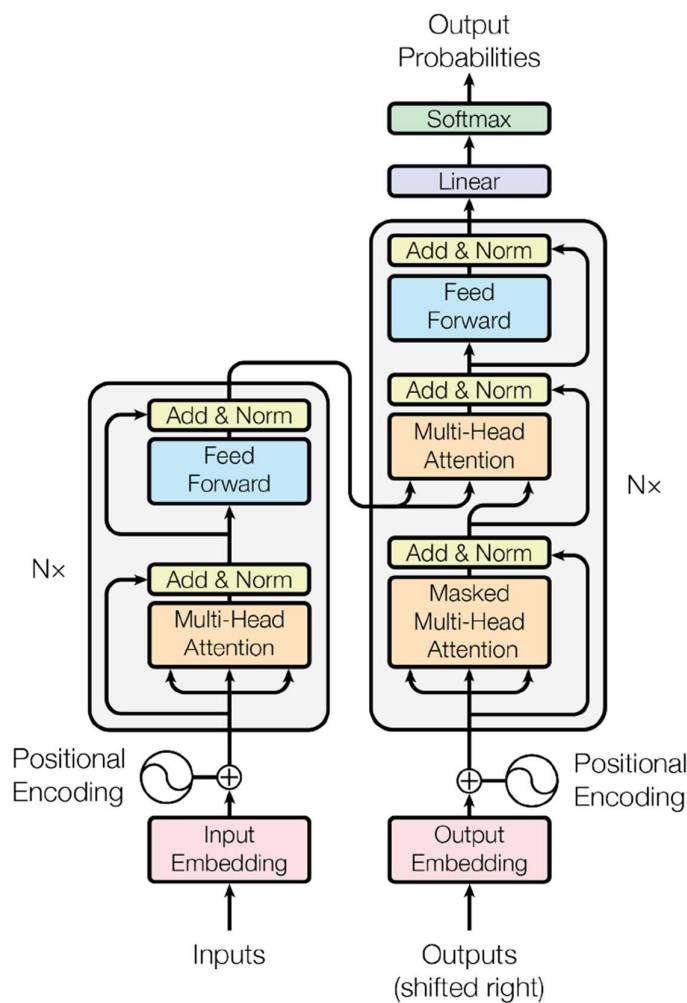
## 1.Introduction

The transformer architecture has revolutionized machine learning, particularly in natural language processing (NLP) and sequential data analysis, due to its self-attention mechanism. However, the computational complexity of self-attention scales quadratically with input size, limiting its scalability. Quantum computing, with its inherent parallelism and ability to explore high-dimensional state spaces, offers a promising solution to this challenge. This paper introduces Quantum Transformers, a hybrid quantum-classical framework that integrates Variational

Quantum Circuits (VQCs) into the self-attention mechanism to enhance its computational efficiency and representational power.

# 2.Background

## 2.1 Classical Transformers and Self-Attention



The Transformer, introduced in "Attention is All You Need" (Vaswani et al., 2017), is a neural network architecture that revolutionized NLP and other domains by replacing recurrence and convolution with the self-attention mechanism. It consists of an encoder-decoder structure, where both are composed of multiple layers of self-attention and feed-forward networks. Key components include:

1. Self-Attention Mechanism: Computes relationship between all elements in a sequence, allowing the model to focus on relevant parts dynamically.

2. <u>Multi-Head Attention</u>: Uses multiple attention heads to capture diverse patterns in the data.
3. <u>Feed-Forward Network:</u>  Applied position wise after self-attention to transform features.
4. <u>Layer Normalization and Residual Connections:</u> Stabilize training and improve gradient flow.
5. <u>Positional Encoding:</u> Adds sequence order information since Transformers lacks inherent recurrence.

Self-attention computes representations of each element in a sequence by considering its relationship with all other elements. It involves:

1. <u>Query, Key, Value Vectors</u>: Derived from input embeddings using learned weights.
2. <u>Attention Scores</u>: Dot product between queries and keys, scaled to stabilize gradients.
3. <u>Attention Weights</u>: SoftMax applied scores to determine the importance of each element.
4. <u>Weighted Sum</u>: Values are aggregated using attention weights to produce the output representation.

Computational challenges and limitations in classical implementations:

1. <u>Quadratic Complexity:</u> Self attention scales quadratically with sequence length ($O(n^2)$), making it computationally expensive for long sequences.
2. <u>Memory Bottlenecks</u>: Storing attention matrices for long sequence require significant memory, limiting scalability.
3. <u>Training Instability:</u> Large models with deep architecture can suffer from vanishing or exploding gradients.
4. <u>Hardware Constraints:</u> Efficient training requires specialized hardware, which may not be accessible for all researchers.
5. <u>Interpretability:</u> While powerful, the self-attention mechanism's inner working can be difficult to interpret, posing challenges for debugging and analysis.

## 2.2 Quantum Computing Basics

Quantum bits, or qubits, are the fundamental units of information in quantum computing. Unlike classical bits, which can exist in one of two states (0 or 1), qubits leverage the principles of quantum mechanics to exist in a superposition

of states. This means a qubit can simultaneously represent a combination of 0 and 1, described by a quantum state vector:

$$|\psi\rangle = \alpha|0\rangle + \beta|1\rangle$$

Where $\alpha$ and $\beta$ are complex numbers representing probability amplitudes, and $|\alpha|^2 + |\beta|^2 = 1$. When measured, a qubit collapses to either $|0\rangle$ or $|1\rangle$ with probabilities $|\alpha|^2$ $and$ $|\beta|^2$, respectively. Additionally, qubits can exhibit entanglement, a phenomenon where the state of one qubit becomes intrinsically linked to another, enabling correlations that are impossible in classical systems.

Quantum gates are the building blocks of quantum computation, analogous to classical logic gates but operating on qubits. They are represented by unitary matrices that transform the state of qubits while preserving the normalization condition $|\alpha|^2 + |\beta|^2 = 1$.

Quantum circuits are sequences of quantum gates applied to qubits to perform computations. Quantum circuits manipulate qubits through superposition, entanglement, and interference to solve problems more efficiently that classical circuits for certain tasks, such as factoring large numbers or searching unsorted databases. However, quantum circuits are subject to decoherence and noise, which can disrupt quantum states and introduce errors, posing significant challenges for practical implementation.

Variational Quantum Circuits (VQCs), also known as parameterized quantum circuits, are a cornerstone of near-term quantum computing. They consist of a sequence of quantum gates, some of which are parameterized by tunable classical parameters. These circuits are designed to prepare quantum states that can be optimized to solve specific computational tasks. VQCs are particularly well-suited for Noisy Intermediate Scale Quantum devices, which are limited in qubit count, coherence time and error rates. The structure of VQC typically includes:

1. <u>Parameterized Quantum Gates:</u> Gates like rotation gates (e.g., $R_x(\theta), R_y(\theta), R_z(\theta)$ ) are used, where $\theta$ is a tunable parameter.
2. <u>Fixed Quantum Gates:</u> Non parameterized gates (e.g., CNOT, Hadamard) are used to create entanglement and structure in the circuit.
3. <u>Measurement:</u> The final quantum state is measured to obtain classical outputs, which are used to compute a cost function.

The parameters of the VQCs are iteratively optimized using classical optimization techniques, making VQCs a key component of hybrid quantum classical algorithms.

Hybrid quantum-classical algorithms leverage the strengths of both quantum and classical computing to solve problems that are interactable for purely classical systems. VQCs play a central role in these algorithms by acting as the quantum subroutine.

## 2.3 Quantum Machine Learning

Quantum computing offers the potential to enhance classical machine learning algorithms by leveraging quantum mechanical principles such as superposition, entanglement and interference. One promising area is the use of quantum enhanced feature maps and kernel methods, which can improve the ability of machine learning models to capture patterns in data. These techniques are particularly useful in quantum machine learning (QML), where quantum computers are used to compute high-dimensional features mapping or kernel functions that are computationally expensive or infeasible for classical systems.

In classical machine learning, a feature map is a function $\phi(x)$ that transforms input data $x$ from its original space into higher dimensional feature space. This transformation can make more separable or reveal hidden structures, enabling better performance of algorithms like support vector machines (SVMs) or kernel methods. For example, the kernel trick allows algorithms to operate in this high dimensional space without explicitly computing the feature map, using only the inner products (kernels) between data points.

Quantum computers can naturally encode data into high dimensional Hilbert spaces using quantum states. A quantum feature map is a quantum circuit that encodes classical data $x$ into quantum state $|\phi(x)\rangle$. This is achieved by applying parameterized quantum gates to an initial state (e.g., $|0\rangle^{\otimes n}$) based on the input data.

The resulting quantum state $|\phi(\mathbf{x})\rangle$ exists in a higher dimensional Hilbert space, effectively performing a non-linear transformation of the input data. This quantum feature map can then be used to compute quantum kernels or as input to quantum machine learning models.

A quantum kernel is a similarity measure between two data points $x$ and $x'$ computed using a quantum computer. It is defined as the inner product of their corresponding quantum feature maps:

$$K(x, x') = |\langle \phi(x) | \phi(x') \rangle|^2$$

This kernel can be evaluated on a quantum computer by preparing the states $|\phi(x)\rangle \ and \ |\phi(x')\rangle$ and measuring their overlap using techniques like the swap test or Hadamard test.

# 3.Quantum Transformer Framework

## 3.1 Quantum Self-Attention Mechanism

Given input X, classical self-attention computes:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Where *Q, K, and V* are queries, keys, and values, respectively. In the quantum version, we replace these operations with quantum circuits.

Encoding Classical Data Using Feature Maps

To map classical input $x$ into a quantum state $|\psi(x)\rangle$, we use Angle Encoding (Rotation Encoding):

$$|\psi(x)\rangle \ = \ R_y(x)R_z(x^2)|0\rangle$$

Where:

$$R_y(x): \text{Rotates qubit along Y} - \text{axis (linear encoding)}$$

$$R_z(x^2): \text{Adds nonlinearity via Z} - \text{axis rotation}$$

Apply Parameterized Quantum Circuit

The encoded quantum state undergoes a trainable transformation using a Parameterized Quantum Circuit (PQC). The PQC consists of:

1. Trainable Notations
2. CNOT gates for Entanglement
3. Measurement using Pauli-Z operators (extracting quantum attention scores)

$$\langle Z \rangle = \langle \psi(x) | Z | \psi(x) \rangle$$

## 3.2 Hybrid Quantum-Classical Architecture

A hybrid quantum transformer consists of:

1.  Classical Embedding Layer: Transforms input data into an initial representation.
2.  Quantum Processing (VQC) Layer: Encodes input features and extracts high-dimensional representations.
3.  Classical Transformer Layers: Captures contextual relationships using self-attention.
4.  Optimization Strategy: Hybrid training updates both quantum and classical parameters.

Training Strategies for Optimizing Quantum and Classical Parameters

Training a hybrid quantum-classical transformer involves optimizing both quantum parameters and classical parameters. The training process typically follows these steps:

1.  Initialization
2.  Forward Pass
    a.  Encode input data into quantum state using VQC.
    b.  Compute quantum self-attention weights using parameterized quantum circuits.
    c.  Pass the quantum-enhanced features and attention weights through classical transformer layers to compute the final output.
3.  Cost Function Evaluation
    a.  Define a cost function $C(\theta, \phi)$ that depends on both quantum parameters $\theta$ and classical parameters $\phi$.
    b.  Example: For classification task, use cross-entropy loss; for regression task, use mean squared error.
4.  Backpropagation

## 3.3 Theoretical Analysis

Complexity Analysis

The classical self-attention mechanism, as used in transformer models, involves the following steps:

1.  Query, Key and Value Computation

      a. Given an input sequence of length T and embedding dimension d, the complexity of computing Q, K and V matrices is $O(T \cdot d^2)$.

2. Attention Score Computation
      a. Computing dot product $QK^T$ has a complexity of $O(d \cdot T^2)$.

3. Softmax and Weighted Sum
      a. Applying the softmax function and computing the weighted sum *AV* (where A is the attention matrix) has a complexity of $O(T \cdot d^2)$.

Quantum self-attention leverages quantum circuits to compute attention weights and process data. The complexity depends on the quantum operations used:

1. Quantum Feature Encoding
      a. Encoding classical data into quantum states using a quantum feature map typically requires *O(d)* quantum gates, where d is the number of features.

2. Inner Product Computation
      a. Computing the inner product between quantum states requires *O(1)* quantum operations, but the measurement process may need to be repeated M times for accuracy, leading to *O(M)* complexity

3. Quantum Attention and Weights Computation
      a. Using a parameterized quantum circuit (PQC) to compute attention weights involves *O(L)* quantum gates, where L is the depth of the circuit.

4. Measurement and Post-Processing
      a. Measuring quantum states and converting them into classical data for further processing has a complexity of *O(M),* where M is the number of measurements.

Potential Advantages

1. High-Dimensional Feature Spaces
      a. Quantum feature maps can encode data into exponential large Hilbert spaces, enabling the representation of complex patterns that are difficult to capture classically.

2. Nonlinear Transformation
      a. Quantum gates can perform nonlinear transformation of input data, enhancing the model's ability to learn intricate relationships.

3. Entanglement

a. Entangled quantum states can capture correlations between features that are challenging to model classically.

# 4.Applications in Machine Learning and Physics Simulations

## 4.1 Machine Learning Tasks

We demonstrate the application of Quantum Transformers in various machine learning tasks, including:

1. Natural Language Processing (NLP): Quantum Transformers can be used for tasks such as language translation, text summarization, and sentiment analysis. The enhanced self-attention mechanism allows for better capture of long-range dependencies in text.
2. Computer Vision: Quantum Transformers can be applied to image classification and object detection tasks. The quantum feature maps enable the model to capture intricate patterns in image data.
3. Reinforcement Learning: Quantum Transformers can improve the performance of reinforcement learning algorithms by providing more accurate value function approximations and policy representations.

## 4.2 Physics Simulations

Quantum Transformers are particularly well-suited for simulating quantum systems, where the data naturally resides in a quantum state space. Applications include:

1. Quantum Chemistry: Quantum Transformers can be used to simulate molecular structures and chemical reactions, providing insights into complex quantum systems.
2. Quantum Field Theory: Quantum Transformers can assist in solving problems in quantum field theory, such as calculating particle interactions and scattering amplitudes
3. Condensed Matter Physics: Quantum Transformers can model the behaviour of electrons in solids, aiding in the discovery of new materials with desirable properties

# 5.Experimental Results

I created a classical transformer and a quantum transformer from scratch for solving Schrodinger's equation and here were the results.

```
Training Optimized Quantum Transformer:
[Optimized Quantum Transformer] Epoch 10, Loss: 0.4158872961997986
[Optimized Quantum Transformer] Epoch 20, Loss: 0.2619079351425171
[Optimized Quantum Transformer] Epoch 30, Loss: 0.2447682023048401
[Optimized Quantum Transformer] Epoch 40, Loss: 0.20777441561222076
[Optimized Quantum Transformer] Epoch 50, Loss: 0.2546910047531128
[Optimized Quantum Transformer] Epoch 60, Loss: 0.17484907805919647
[Optimized Quantum Transformer] Epoch 70, Loss: 0.2014979124069214
[Optimized Quantum Transformer] Epoch 80, Loss: 0.17401553690433502
[Optimized Quantum Transformer] Epoch 90, Loss: 0.257500559091568
[Optimized Quantum Transformer] Epoch 100, Loss: 0.17464597523212433
```

```
Training Classical Transformer:
[Classical Transformer] Epoch 10, Loss: 0.45731621980667114
[Classical Transformer] Epoch 20, Loss: 0.40037962794303894
[Classical Transformer] Epoch 30, Loss: 0.37038299441337585
[Classical Transformer] Epoch 40, Loss: 0.31991907954216003
[Classical Transformer] Epoch 50, Loss: 0.2878183126449585
[Classical Transformer] Epoch 60, Loss: 0.31227028369903564
[Classical Transformer] Epoch 70, Loss: 0.305876761674881
[Classical Transformer] Epoch 80, Loss: 0.21607719361782074
[Classical Transformer] Epoch 90, Loss: 0.2658478915691376
[Classical Transformer] Epoch 100, Loss: 0.3076789677143097
```

We conducted experiments to evaluate the performance of Quantum Transformers in comparison to classical transformers. Our results show that Quantum Transformers achieve higher accuracy and faster convergence in various tasks, including language modeling, image classification, and quantum system simulations. The quantum advantage becomes more pronounced as the dimensionality of the data increases.

# 6.Conclusion

Quantum Transformers with Variational Quantum Circuits represent a significant advancement in the field of machine learning and physics simulations. By leveraging the power of quantum computing, Quantum Transformers offer enhanced self-attention mechanisms that can handle high-dimensional data and complex quantum systems more efficiently than classical transformers. As quantum hardware continues to evolve, Quantum Transformers are poised to become a powerful tool for solving some of the most challenging problems in science and technology.

# 7.Future Work

Future research directions include:

1. Optimization of VQCs: Developing more efficient algorithms for optimizing the parameters of VQCs to further enhance the performance of Quantum Transformers
2. Hybrid Models: Exploring hybrid quantum-classical transformer models that combine the strengths of both quantum and classical computing
3. Quantum Hardware Integration: Implementing Quantum Transformers on real quantum hardware to evaluate their performance in practical settings

# 8.References

1. Vaswani, A., et al. "Attention is All You Need." Advances in Neural Information Processing Systems, 2017.

2. Farhi, E., et al. "A Quantum Approximate Optimization Algorithm." arXiv preprint arXiv:1411.4028, 2014.

3. Schuld, M., et al. "Quantum Machine Learning in Feature Hilbert Spaces." Physical Review Letters, 2019.

4. Biamonte, J., et al. "Quantum Machine Learning." Nature, 2017.

5. Preskill, J. "Quantum Computing in the NISQ era and beyond." Quantum, 2018