

# FIGHTER: UNVEILING THE GRAPH CONVOLUTIONAL NATURE OF TRANSFORMERS IN TIME SERIES MODELING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Transformers have achieved remarkable success in time series modeling, yet their internal mechanisms remain opaque. This work demystifies the Transformer encoder by establishing its fundamental equivalence to a Graph Convolutional Network (GCN). We show that in the forward pass, the attention distribution matrix serves as a dynamic adjacency matrix, and its composition with subsequent transformations performs computations analogous to graph convolution. Moreover, we demonstrate that in the backward pass, the update dynamics of value and feed-forward projections mirror those of GCN parameters. Building on this unified theoretical reinterpretation, we propose **Fighter** (Flexible Graph Convolutional Transformer), a streamlined architecture that removes redundant linear projections and incorporates multi-hop graph aggregation. This perspective yields an explicit and interpretable representation of temporal dependencies across different scales, naturally expressed as graph edges. Experiments on standard forecasting benchmarks confirm that Fighter achieves competitive performance while providing clearer mechanistic interpretability of its predictions.

## 1 INTRODUCTION

Time series data is prevalent in the current data-centric world. Effective modeling of time series is essential for informed decision-making and strategic planning across various domains, including finance (Tsay, 2005), transportation (Vlahogianni & Karlaftis, 2013), and energy management (Wang et al., 2007). The success of Transformer models (Vaswani et al., 2017) in natural language processing has quickly drawn the attention of researchers working on time series (Wen et al., 2023). Thanks to their ability to capture long-range dependencies, Transformer-based models have emerged as promising solutions for overcoming common challenges in time series modeling.

However, the internal workings of Transformers remain poorly understood (Rudin, 2019; Zeng et al., 2023) in time series, highlighting the need for a deeper investigation:

- **How does the Transformer process input time series data?** Specifically, what role do attention mechanisms and feed-forward networks (FFNs) play?
- **What are the update dynamics of Transformers?** In particular, how are the parameters in the attention mechanisms and FFNs updated?

Understanding these questions is crucial for the reliable application of Transformers in time series data, such as finance (Masini et al., 2023) and healthcare (Wickstrøm et al., 2020; Morid et al., 2023), and will offer insights into the design of more effective model architectures.

To this end, this work presents a unified theoretical reinterpretation that establishes a fundamental equivalence between the Transformer encoder and Graph Convolutional Networks (GCNs) (Kipf & Welling, 2017). GCNs are powerful architectures for modeling graphs, defined by a static adjacency matrix and feature matrix that represent fixed edge connections and node attributes, respectively (Hamilton et al., 2017; Chami et al., 2022). A GCN operates similarly to a multilayer perceptron (MLP), with its key advancement lying in aggregating features across multiple hops based on the adjacency matrix at the start of each layer (Wu et al., 2019; Zhang et al., 2025). We demonstrate that, in the forward pass, the Transformer attention distribution matrix (*i.e.*,  $\text{softmax}(\mathbf{Q}\mathbf{K}^\top/\sqrt{d} + M)$ ),

where  $Q$ ,  $K$ , and  $M$  denote query, key, and mask matrices) serves as a dynamic counterpart to the adjacency matrix for feature aggregation. Unlike the static adjacency matrix used in GCNs, which corresponds to fixed graph connections, this distribution matrix is learnable via the parameters of the query and key matrices. The subsequent linear projections in the value matrix and feed-forward network (FFN)<sup>1</sup> perform computations that resemble MLP-based feature extraction in GCNs. During training, we show that the updates to the value matrix and FFN parameters mirror those in GCNs, focusing on feature extraction. Meanwhile, the parameters of  $Q$  and  $K$  facilitate an adaptive update, which is responsible for a dynamic and learnable adjacency matrix, in contrast to the static connections involved in GCNs.

Building on this unified perspective, we propose **Fighter** (Flexible Graph Convolutional Transformer), a streamlined architecture that removes redundant linear projections and integrates multi-hop graph aggregation for an explicit and interpretable representation of temporal dependencies across various scales. Lastly, experiments on standard forecasting benchmarks confirm that Fighter achieves competitive performance while offering clearer mechanistic interpretability of its predictions. Our key contributions are listed as follows:

- We establish a fundamental theoretical equivalence between the Transformer encoder and Graph Convolutional Networks (GCNs), revealing that the Transformer attention distribution matrix acts as a dynamic adjacency matrix for feature aggregation, while parameter updates to the value matrix and FFN during training mirror GCN dynamics. This unified framework demystifies the internal workings of Transformer encoders in time series modeling, bridging sequence and graph paradigms.
- We propose Fighter (Flexible Graph Convolutional Transformer), a novel architecture that leverages this equivalence by eliminating redundant linear projections and integrating multi-hop graph aggregation, enabling explicit and interpretable representations of multi-scale temporal dependencies as graph edges.
- We empirically evaluate Fighter on standard time series forecasting benchmarks, demonstrating its ability to effectively model temporal dependencies while providing enhanced interpretability through graph-based insights, as evidenced by detailed analyses of its attention distribution matrix.

## 2 RELATED WORKS

**Transformer-Based Time Series Modeling.** Transformer-based models have significantly advanced time series modeling by capturing long-range dependencies through sophisticated attention mechanisms (Vaswani et al., 2017). Recent efforts have explored diverse improvements in time series forecasting, particularly in computational efficiency and temporal periodicity modeling. For efficiency, Informer introduced ProbSparse attention to process long sequences with reduced computational demands (Zhou et al., 2021), followed by Pyraformer, which employs pyramidal attention to capture multi-scale dependencies efficiently (Liu et al., 2022). Triformer further optimizes memory usage with triangle attention (Cirstea et al., 2022), and PatchTST leverages patch-based processing to minimize resource requirements (Nie et al., 2023). To capture periodicity, Autoformer integrates auto-correlation to model recurring patterns effectively (Wu et al., 2021), while FEDformer enhances accuracy by combining frequency-domain attention with decomposition (Zhou et al., 2022). Despite these advancements, the opaque internal workings of Transformers and redundant linear projections often limit their interpretability and efficiency (Zeng et al., 2023). More recently, PFFormer introduces position-free embeddings to better capture complex dependencies in extreme adaptive multivariate settings (Li & Anastasiu, 2025), and PENGUIN leverages periodic-nested group attention with relative bias to explicitly model multiple coexisting temporal periodicities (Sun et al., 2025). This work addresses these limitations by reinterpreting Transformers through a graph-based lens, enabling clearer mechanistic insights and streamlined computations.

**Graph-Based Methods and Interpretability.** Graph Convolutional Networks (GCNs) effectively capture graph-structured data by aggregating features across multi-hop neighborhoods using static adjacency matrices (Kipf & Welling, 2017; Wu et al., 2019; Yu et al., 2018). In multivariate time series forecasting, models such as GTS (Shang et al., 2021), MTGNN (Wu et al., 2020), AGCRN (Bai et al., 2020), and StemGNN (Cao et al., 2020) typically construct (static or adaptive) graphs over the channel dimension and handle temporal dynamics separately via convolutions or recurrence,

<sup>1</sup>A feed-forward network is a 2-layer multilayer perceptron.

108 leaving the temporal axis itself unmodeled as a graph. FourierGNN (Yi et al., 2023) stands out as a  
 109 rare exception by constructing a joint space-time graph over both channels and time steps, yet its fully  
 110 connected design incurs prohibitive complexity, necessitating sliding windows that severely constrain  
 111 long-range temporal modeling. Some studies have begun treating sequences as graphs (Li et al.,  
 112 2024; Joshi, 2025), yet they lack a systematic theoretical comparison of Transformers and GCNs  
 113 across both forward and backward passes, fail to eliminate redundant linear projections, and do not  
 114 effectively leverage explicit multi-hop aggregation. Existing interpretability efforts in time series,  
 115 such as attention visualization (Vaswani et al., 2017) or feature attribution methods (Wickstrøm et al.,  
 116 2020), are largely post-hoc and fail to provide intrinsic mechanistic insights. In contrast, our Fighter  
 117 architecture bridges Transformers and GCNs by casting the attention distribution matrix as a dynamic,  
 118 learnable adjacency matrix defined exclusively over the temporal dimension, removes redundant  
 119 linear projections for streamlined computation, and incorporates multi-hop aggregation to explicitly  
 120 represent temporal dependencies as interpretable graph edges, significantly enhancing both efficiency  
 121 and interpretability.

### 122 3 BACKGROUND

123 **Notation.**<sup>2</sup> Consider a graph  $\mathcal{G} = (\mathbf{X}_{n \times d}, \mathbf{A}_{n \times n}) \in \mathcal{G}$ , where  $\mathbf{X}_{n \times d}$  is an  $n \times d$  feature matrix  
 124 comprising all node feature vectors  $\mathbf{x}_i \in \mathbb{R}^d$ , for  $i \in \mathbb{N}_n$  with  $\mathbb{N}_n := \{1, \dots, n\}$ , and  $\mathbf{A}_{n \times n}$  is the  
 125 adjacency matrix encoding the graph’s edge connections. Additionally, let  $\mathbf{S}_{1:S} = (\mathbf{x}_1, \dots, \mathbf{x}_S) \in \mathcal{S}$   
 126 denote a sequence of length  $S$ , where each  $\mathbf{x}_s \in \mathbb{R}^d$  is the feature vector for the  $s$ -th element,  $s \in \mathbb{N}_S$ ,  
 127 and the collection of these vectors forms an  $S \times d$  feature matrix, denoted  $\mathbf{X}_{S \times d}$ . In both graphs and  
 128 sequences, feature vectors  $\mathbf{x}$  are row vectors, written as  $[x_j]_d^\top = [x_1, \dots, x_d]$ . The  $i$ -th row and  $j$ -th  
 129 column of a feature matrix, corresponding to the  $i$ -th node or element and  $j$ -th feature, are denoted  
 130 by  $\mathbf{X}_{(i,:)} := \mathbf{e}_i^\top \mathbf{X}$  and  $\mathbf{X}_{(:,j)} := \mathbf{X} \mathbf{e}_j$ , respectively, where  $\mathbf{e}_i$  is a standard basis vector with a 1 in  
 131 the  $i$ -th position and 0 elsewhere. The identity matrix is denoted by  $\mathbf{I}$ . A block diagonal matrix with  
 132 entry matrices  $\mathbf{M}_1, \dots, \mathbf{M}_m$  is denoted  $\text{blkdiag}(\mathbf{M}_1, \dots, \mathbf{M}_m)$ , and if all  $m$  entries are identical, it  
 133 is simplified to  $\text{blkdiag}(\mathbf{M}; m)$ .

134 **Flexible graph convolutional network (GCN)** extends traditional GCN by enabling adaptable  
 135 weight assignments for aggregated features (Krishnagopal & Ruiz, 2023). Unlike conventional GCN,  
 136 which uses  $(\mathbf{A} + \mathbf{I})^\kappa \mathbf{X}$  for feature aggregation—limiting flexibility by applying uniform weights  
 137 to features from different hops (*i.e.*, neighbor distances) within a single layer  $\sigma((\mathbf{A} + \mathbf{I})^\kappa \mathbf{X} \mathbf{W})$   
 138 with weight matrix  $\mathbf{W}$ —a flexible GCN  $\mathbf{X}_{\text{GCN}}^{(L)}$  unfolds these features and assigns distinct weights  
 139 (Abu-El-Haija et al., 2019; Zhang et al., 2025). This can be expressed as:

$$140 \mathbf{X}_{\text{GCN}}^{(\ell)} = \sigma \left( \mathbf{A}^{[\kappa_\ell]} \text{blkdiag}(\mathbf{X}_{\text{GCN}}^{(\ell-1)}; \kappa_\ell) \cdot \mathbf{W}^{(\ell)} \right), \ell \in \mathbb{N}_{L-1}$$

$$141 \mathbf{X}_{\text{GCN}}^{(L)} = \mathbf{A}^{[\kappa_L]} \text{blkdiag}(\mathbf{X}_{\text{GCN}}^{(L-1)}; \kappa_L) \cdot \mathbf{W}^{(L)}, \quad (1)$$

142 where  $\mathbf{A}^{[\kappa]} := \|\|_{i=0}^{\kappa-1} \mathbf{A}^i = [\mathbf{I}, \mathbf{A}, \dots, \mathbf{A}^{\kappa-1}]$  is an  $n \times \kappa n$  matrix formed by concatenation  $\|\|$ ,  $\mathbf{W}^{(\ell)}$   
 143 is the weight matrix at layer  $\ell$  with dimensions  $h_{\ell-1} \times h_\ell$ ,  $h_\ell$  represents the width of layer  $\ell$  with  
 144  $h_0 = d$ , and  $\kappa_\ell$  denotes the hop distance at the  $\ell$ -th layer.  $\mathbf{X}_{\text{GCN}}^{(0)}$  is the input feature matrix, and  $\sigma$   
 145 denotes the activation function (*e.g.*, ReLU). Such a flexible GCN resembles an MLP but aggregates  
 146 features across multiple hops at the start of each layer. By restricting the multi-hop aggregation,  
 147 governed by  $\mathbf{A}^{[\kappa_\ell]} \text{blkdiag}(\mathbf{X}_{\text{GCN}}^{(\ell-1)}; \kappa_\ell)$ , to single-hop aggregation (*i.e.*, using only the  $\mathbf{A}^1$  term), the  
 148 feature representation  $\mathbf{X}_{\text{GCN}}^{(\ell)}$  at layer  $\ell \in \mathbb{N}_{L-1}$  is simplified as:

$$149 \mathbf{X}_{\text{GCN}}^{(\ell)} = \sigma \left( \mathbf{A} \mathbf{X}_{\text{GCN}}^{(\ell-1)} \mathbf{W}^{(\ell)} \right), \ell \in \mathbb{N}_{L-1}. \quad (2)$$

150 **Transformer**, originally developed for natural language processing, comprises two main components:  
 151 an encoder and a decoder (Vaswani et al., 2017). For representation learning in time series modeling,  
 152 the Transformer encoder  $\mathbf{X}_{\text{TF}}^{(L)}$ , with its attention layer  $\mathbf{X}_{\text{Att}}^{(\ell)}$  and FFN layer  $\mathbf{X}_{\text{FFN}}^{(\ell)}$ , has demonstrated  
 153 high effectiveness (Zerveas et al., 2021; Yıldız et al., 2022). Specifically, the Transformer encoder

154 <sup>2</sup>See the notation table in Appendix A.1.

with self-attention can be formulated as:

$$\begin{aligned} \mathbf{X}_{\text{Att}}^{(\ell)} &= \text{smx} \left( d^{-1/2} \mathbf{X}_{\text{FFN}}^{(\ell-1)} \mathbf{W}_Q^{(\ell)} \left( \mathbf{X}_{\text{FFN}}^{(\ell-1)} \mathbf{W}_K^{(\ell)} \right)^\top \right) \mathbf{X}_{\text{FFN}}^{(\ell-1)} \mathbf{W}_V^{(\ell)}, \ell \in \mathbb{N}_{L-1} \\ \mathbf{X}_{\text{FFN}}^{(\ell)} &= \sigma \left( \mathbf{X}_{\text{Att}}^{(\ell)} \mathbf{W}_{\text{FFN}}^{(\ell,1)} + \mathbf{b}_{\text{FFN}}^{(\ell,1)} \right) \mathbf{W}_{\text{FFN}}^{(\ell,2)} + \mathbf{b}_{\text{FFN}}^{(\ell,2)} \\ \mathbf{X}_{\text{TF}}^{(L)} &= \mathbf{X}_{\text{FFN}}^{(L-1)} \mathbf{W}^{(L)}, \end{aligned} \quad (3)$$

where  $\text{smx}(\cdot)$  denotes a row-wise softmax operation;  $\mathbf{W}_Q^{(\ell)}$ ,  $\mathbf{W}_K^{(\ell)}$ , and  $\mathbf{W}_V^{(\ell)}$  are the query, key, and value weight matrices at layer  $\ell$ , with dimensions  $h_{\ell-1} \times p$ ,  $h_{\ell-1} \times p$ ,  $h_{\ell-1} \times v$  respectively;  $\mathbf{W}_{\text{FFN}}^{(\ell,1)}$  and  $\mathbf{W}_{\text{FFN}}^{(\ell,2)}$  are the weight matrices for the first and second sublayers of the FFN at layer  $\ell$ , with dimension  $v \times h_\ell$  and  $h_\ell \times h_\ell$ ;  $\mathbf{b}_{\text{FFN}}^{(\ell,\cdot)}$  represents the bias terms;  $\mathbf{X}_{\text{FFN}}^{(0)}$  is the input feature matrix of size  $S \times d$  with  $d = h_0$ ; and  $\sigma$  is the row-wise activation function (*e.g.*, ReLU). For clarity and conciseness, this paper discusses a single-head Transformer encoder, omitting tokenization, positional encodings, layer normalization, and residual connections from the formulation. These components are fully incorporated in our experimental implementation.

## 4 TRANSFORMERS AS DYNAMIC GRAPH CONVOLUTION FOR TIME SERIES

We begin by analyzing the forward pass of Transformer encoders, where the attention distribution matrix serves as a dynamic adjacency matrix, and its composition with subsequent transformations performs graph convolution-like computations. Extending this analysis to the training phase, we show that gradients update the value matrix and FFN parameters similarly to GCNs, while the query and key weight matrices enable adaptive updates to a learnable adjacency matrix. Building on these findings, we introduce Fighter, a flexible graph convolutional transformer that removes redundant linear projections and incorporates multi-hop graph aggregation for interpretable modeling of temporal dependencies.

### 4.1 FORWARD PASS: THE ATTENTION DISTRIBUTION MATRIX AS A DYNAMIC ADJACENCY

To reveal the graph convolutional nature of Transformers, we reformulate the Transformer encoder and compare it to GCNs. For a Transformer encoder as expressed in Equation 3, the learned hidden feature  $\mathbf{X}_{\text{TF}}^{(\ell)}$  at layer  $\ell \in \mathbb{N}_{L-1}$  can be reorganized, omitting bias terms for simplicity and defining  $\mathbf{W}_V^{(\ell,1)} := \mathbf{W}_V^{(\ell)} \mathbf{W}_{\text{FFN}}^{(\ell,1)}$ , as:

$$\mathbf{X}_{\text{TF}}^{(\ell)} = \sigma \left( \text{smx} \left( d^{-1/2} \mathbf{X}_{\text{TF}}^{(\ell-1)} \mathbf{W}_Q^{(\ell)} \left( \mathbf{X}_{\text{TF}}^{(\ell-1)} \mathbf{W}_K^{(\ell)} \right)^\top \right) \mathbf{X}_{\text{TF}}^{(\ell-1)} \mathbf{W}_V^{(\ell,1)} \right) \cdot \mathbf{W}_{\text{FFN}}^{(\ell,2)}. \quad (4)$$

The term  $\text{smx} \left( d^{-1/2} \mathbf{X}_{\text{TF}}^{(\ell-1)} \mathbf{W}_Q^{(\ell)} \left( \mathbf{X}_{\text{TF}}^{(\ell-1)} \mathbf{W}_K^{(\ell)} \right)^\top \right)$  forms an  $S \times S$  attention distribution matrix, which performs single-hop feature aggregation over the input sequence, analogous to the adjacency matrix  $\mathbf{A}$  in GCNs (Veličković et al., 2018). The derivation of this reformulation is provided in Appendix A.2.

By comparing Equations 4 and 2, we observe that the attention distribution matrix in the Transformer serves as a dynamic analogue to the static adjacency matrix  $\mathbf{A}$  in GCNs, while  $\mathbf{W}_V^{(\ell,1)}$  aligns with  $\mathbf{W}^{(\ell)}$  for feature extraction. That is to say, a Transformer encoder, to some extent, is a single-hop GCN with dynamic adjacency matrix. From this GCN perspective, the additional linear projection  $\mathbf{W}_{\text{FFN}}^{(\ell,2)}$  in the Transformer appears redundant, as GCNs do not employ a post-activation linear projection. Indeed, its role in transforming aggregated features is effectively subsumed by the subsequent layer’s projections, namely  $\mathbf{W}_Q^{(\ell+1)}$ ,  $\mathbf{W}_K^{(\ell+1)}$ , and  $\mathbf{W}_V^{(\ell+1,1)}$ , which collectively manage feature transformations in the next iteration of the recursive architecture. This insight highlights a structural parallel between Transformers and GCNs that has been underexplored. In contrast, prior works (Veličković et al., 2018; Li et al., 2024; Joshi, 2025) focus predominantly on the attention mechanism’s role in isolation, often missing its broader graph convolutional interpretation.

For the final layer,  $\ell = L$ , the Transformer encoder omits the attention mechanism, yielding  $\mathbf{X}_{\text{TF}}^{(L)} = \mathbf{X}_{\text{TF}}^{(L-1)} \mathbf{W}^{(L)}$ . This directly corresponds to the GCN formulation when  $\kappa_L = 1$ , *i.e.*,  $\mathbf{X}_{\text{GCN}}^{(L)} =$

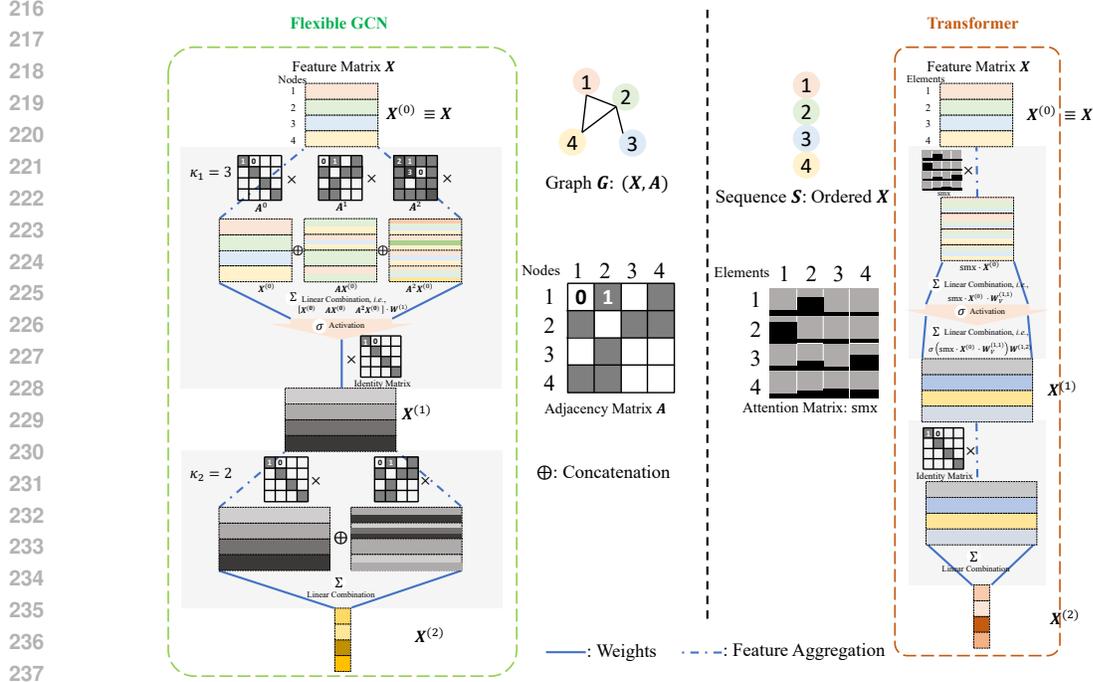


Figure 1: Visual comparison of a Transformer encoder and a flexible GCN, illustrating the equivalence between the Transformer’s attention distribution matrix and the GCN’s adjacency matrix in performing feature aggregation during the forward pass.

$\mathbf{A}^0 \mathbf{X}_{\text{GCN}}^{(L-1)} \mathbf{W}^{(L)}$  (Li et al., 2024). This equivalence underscores that Transformers inherently perform graph-like convolution during the forward pass, with the attention distribution matrix dynamically adapting to the input sequence. A visual comparison of the Transformer encoder and flexible GCN is provided in Figure 1.

**Remark 1.** *The equivalence between Transformer encoders and GCNs stems from the attention distribution matrix dynamically mimicking the GCN’s adjacency matrix for feature aggregation. This equivalence is independent of the specific attention mechanism employed, as the operation within smx can be tailored to the chosen mechanism (e.g., scaled dot-product or variants). This flexibility ensures the graph convolutional interpretation remains applicable across diverse Transformer architectures, unifying sequence and graph-based processing under a common framework.*

## 4.2 BACKWARD PASS: GRADIENT-BASED LEARNING OF DYNAMIC STRUCTURE

Following the graph convolutional interpretation of Transformer encoders in the above Section, we now explore the backward pass to compare the gradient-based learning dynamics of Transformers encoders and GCNs (Gelfand et al., 2000; Ruder, 2016; Chami et al., 2022). We demonstrate that the feature transformation parameters in Transformers, excluding the query and key weight matrices  $\mathbf{W}_Q$  and  $\mathbf{W}_K$ , are updated similarly to those in GCNs. Additionally, we show how the attention distribution matrix, analogous to the GCN’s adjacency matrix, is dynamically learned through the gradients of  $\mathbf{W}_Q$  and  $\mathbf{W}_K$ . For clarity, we consider a two-layer, scalar-valued Transformer encoder and a GCN, as illustrated in Figure 1.

For the Transformer encoder, we simplify Equation 4 by omitting the redundant linear projection  $\mathbf{W}_{\text{FFN}}^{(\ell,2)}$  and setting  $L = 2$ , yielding the output at the final layer:

$$\mathbf{X}_{\text{TF}}^{(2)} = \sigma \left( \text{smx} \left( d^{-1/2} \mathbf{X}^{(0)} \mathbf{W}_Q^{(1)} \left( \mathbf{X}^{(0)} \mathbf{W}_K^{(1)} \right)^\top \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \right) \cdot \mathbf{W}^{(2)}. \quad (5)$$

The gradients with respect to the feature transformation parameters  $\mathbf{W}^{(2)}$  and  $\mathbf{W}_V^{(1,1)}$  ( $i \in \mathbb{N}_{h_1}$ ) are:

$$\begin{aligned} \frac{\partial \mathbf{X}_{\text{TF}}^{(2)}}{\partial \mathbf{W}^{(2)}} &= \sigma \left( \text{smx} \left( d^{-1/2} \mathbf{X}^{(0)} \mathbf{W}_Q^{(1)} \left( \mathbf{X}^{(0)} \mathbf{W}_K^{(1)} \right)^\top \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \right), \\ \frac{\partial \mathbf{X}_{\text{TF}}^{(2)}}{\partial \mathbf{W}_V^{(1,1)}} &= \dot{\sigma} \cdot \text{smx} \left( d^{-1/2} \mathbf{X}^{(0)} \mathbf{W}_Q^{(1)} \left( \mathbf{X}^{(0)} \mathbf{W}_K^{(1)} \right)^\top \right) \mathbf{X}^{(0)} \cdot \mathbf{W}_{(i,:)}^{(2)}, \end{aligned} \quad (6)$$

where  $\dot{\sigma}$  is the derivative of the activation function.

For the GCN, we consider single-hop aggregation in the first layer and set  $\kappa_2 = 1$  in the second layer, yielding the gradients for Equation 2:

$$\frac{\partial \mathbf{X}_{\text{GCN}}^{(2)}}{\partial \mathbf{W}^{(2)}} = \sigma(\mathbf{A} \mathbf{X}^{(0)} \mathbf{W}^{(1)}), \quad \frac{\partial \mathbf{X}_{\text{GCN}}^{(2)}}{\partial \mathbf{W}_{(i,:)}^{(1)}} = \dot{\sigma} \mathbf{A} \mathbf{X}^{(0)} \mathbf{W}_{(i,:)}^{(2)}, \quad (7)$$

which correspond to a specific instance of the multi-hop gradient, with details provided in Appendix A.3.

Comparing Equations 6 and 7, the gradients for feature transformation parameters in both models share a similar structure, with the Transformer’s attention distribution matrix  $\text{smx} \left( d^{-1/2} \mathbf{X}^{(0)} \mathbf{W}_Q^{(1)} \left( \mathbf{X}^{(0)} \mathbf{W}_K^{(1)} \right)^\top \right)$  acting as a dynamic counterpart to the GCN’s static adjacency matrix  $\mathbf{A}$ . This similarity indicates that the Transformer’s feature transformation parameters are updated in a manner analogous to those in GCNs.

However, the Transformer’s attention distribution matrix is dynamically learned through the gradients of  $\mathbf{W}_Q^{(1)}$  and  $\mathbf{W}_K^{(1)}$  ( $i \in \mathbb{N}_p$ ):

$$\frac{\partial \mathbf{X}_{\text{TF}}^{(2)}}{\partial \mathbf{W}_Q^{(1)}} = \left[ \dot{\sigma}_j / \sqrt{d} \mathbf{X}_{(j,:)}^{(0)} \cdot \underbrace{\left( \mathcal{K}_{(i,:)}^{(1)} \right)^\top}_{1 \times d} \cdot \underbrace{\text{blkdiag} \left( \text{smx} \left( \xi_{(j,i;:,i)}^{(1)} \right) \right)}_{S \times S} \cdot \underbrace{\mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)}}_{S \times d \times d \times h_1 \times h_1 \times 1} - \underbrace{\mathcal{K}_{(i,:)}^{(1)}}_{1 \times S} \cdot \underbrace{\left( \text{smx} \left( \xi_{(j,i;:,i)}^{(1)} \right) \right)^\top}_{S \times S} \cdot \underbrace{\text{smx} \left( \xi_{(j,i;:,i)}^{(1)} \right)}_{S \times d} \cdot \underbrace{\mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)}}_{h_1 \times 1} \right]_{S \times d} \quad (8)$$

$$\frac{\partial \mathbf{X}_{\text{TF}}^{(2)}}{\partial \mathbf{W}_K^{(1)}} = \left[ \dot{\sigma}_j / \sqrt{d} \mathbf{X}_{(j,:)}^{(0)} \cdot \underbrace{\left( \mathcal{Q}_{(i,:)}^{(1)} \right)^\top}_{1 \times d} \cdot \underbrace{\text{blkdiag} \left( \text{smx} \left( \xi_{(j,i;:,i)}^{(1)} \right) \right)}_{S \times S} \cdot \underbrace{\mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)}}_{S \times d \times d \times h_1 \times h_1 \times 1} - \underbrace{\mathcal{Q}_{(i,:)}^{(1)} \left( \text{smx} \left( \xi_{(j,i;:,i)}^{(1)} \right) \right)^\top}_{1 \times 1} \cdot \underbrace{\text{smx} \left( \xi_{(j,i;:,i)}^{(1)} \right)}_{S \times S} \cdot \underbrace{\mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)}}_{S \times d \times d \times h_1 \times h_1 \times 1} \right]_{S \times d} \quad (9)$$

where  $\mathcal{Q} := \mathbf{X}^{(0)} \mathbf{W}_Q^{(1)}$ ,  $\mathcal{K} := \mathbf{X}^{(0)} \mathbf{W}_K^{(1)}$ ,  $\xi_{(j,i;:,i)}^{(1)} := \mathcal{Q}_{(j,i)}^{(1)} \mathcal{K}_{(i,:)}^{(1)\top} / \sqrt{d}$ , and the derivation is provided in Appendix A.4. Focusing on the query gradient (Equation 8), the update depends on both the input features  $\mathbf{X}_{(j,:)}^{(0)}$  and a scalar  $\omega_j$ , which assigns varying importance to each sequence element during training. This dynamic adaptation of the attention distribution matrix, absent in the fixed adjacency matrix of GCNs, interprets the Transformer’s expressiveness.

### 4.3 FIGHTER: EFFICIENT MULTI-HOP GRAPH CONVOLUTIONAL TRANSFORMER

Building on the unified theoretical reinterpretation of Transformer encoders as GCNs, we propose **Fighter** (Flexible Graph Convolutional Transformer), a streamlined architecture that eliminates redundant linear projections and incorporates multi-hop graph aggregation. This design enhances expressive power by extending the dynamic adjacency matrix’s capabilities in the forward pass while reducing computational overhead. Specifically, the forward pass of Fighter for the learned hidden features  $\mathbf{X}_{\text{FT}}^{(\ell)}$  at layer  $\ell \in \mathbb{N}_{L-1}$  and the final output  $\mathbf{X}_{\text{FT}}^{(L)}$  is expressed as:

$$\begin{aligned} \mathbf{X}_{\text{FT}}^{(\ell)} &= \sigma \left( \text{smx}^{[\kappa_\ell]} \left( d^{-1/2} \mathbf{X}_{\text{FT}}^{(\ell-1)} \mathbf{W}_Q^{(\ell)} \left( \mathbf{X}_{\text{FT}}^{(\ell-1)} \mathbf{W}_K^{(\ell)} \right)^\top \right) \text{blkdiag} \left( \mathbf{X}_{\text{FT}}^{(\ell-1)}; \kappa_\ell \right) \cdot \mathbf{W}^{(\ell)} \right) \\ \mathbf{X}_{\text{FT}}^{(L)} &= \text{smx}^{[\kappa_L]} \left( d^{-1/2} \mathbf{X}_{\text{FT}}^{(L-1)} \mathbf{W}_Q^{(L)} \left( \mathbf{X}_{\text{FT}}^{(L-1)} \mathbf{W}_K^{(L)} \right)^\top \right) \text{blkdiag} \left( \mathbf{X}_{\text{FT}}^{(L-1)}; \kappa_L \right) \cdot \mathbf{W}^{(L)} \end{aligned} \quad (10)$$

where the value and feed-forward projection is **streamlined**, and **multi-hop aggregation** is enabled through the raised attention distribution matrix  $\text{smx}^{[\kappa_\ell]}$  and the block-diagonal replication  $\text{blkdiag}(\cdot; \kappa_\ell)$ . A visualization of Fighter is provided in Figure 2.

324 Fighter reinterprets the Transformer encoder through the lens of a GCN with a dynamic adjacency  
 325 matrix, as established in the forward pass analysis (Section 4.1). Here, the attention distribution  
 326 matrix  $\text{smx} \left( d^{-1/2} \mathbf{X}_{\text{TF}}^{(\ell-1)} \mathbf{W}_Q^{(\ell)} \left( \mathbf{X}_{\text{TF}}^{(\ell-1)} \mathbf{W}_K^{(\ell)} \right)^\top \right)$  functions as a learnable counterpart to the static  
 327 adjacency matrix  $\mathbf{A}$  in GCNs. To mitigate the redundancy of the additional linear projection  $\mathbf{W}_{\text{FFN}}^{(\ell,2)}$ —  
 328 which, as shown, is effectively subsumed by subsequent layer projections—Fighter removes it entirely.  
 329 Instead, it relies on the combined value and initial FFN projection (consolidated into  $\mathbf{W}^{(\ell)}$ ) alongside  
 330 transformations in ensuing layers for efficient feature extraction, thereby simplifying the architecture  
 331 without sacrificing performance.  
 332

333 Furthermore, to overcome the limita-  
 334 tions of single-hop aggregation in  
 335 standard Transformers (Equation 4),  
 336 Fighter integrates multi-hop aggrega-  
 337 tion, drawing inspiration from the  
 338 flexible GCN (Equation 1). By rais-  
 339 ing the attention distribution matrix  
 340 to powers up to  $\kappa_\ell - 1$  and replicat-  
 341 ing input features accordingly, Fighter  
 342 explicitly captures multi-scale tempo-  
 343 ral dependencies. This models indi-  
 344 rect relationships across time steps  
 345 as multi-hop graph connections, en-  
 346 abling Fighter to better handle com-  
 347 plex, long-range patterns inherent in  
 time series data.

348 In the backward pass, Fighter retains  
 349 the dynamic learning of the adjacency  
 350 matrix through gradients with respect  
 351 to  $\mathbf{W}_Q^{(\ell)}$  and  $\mathbf{W}_K^{(\ell)}$ , as analyzed in  
 352 Section 4.2. These updates adapt the  
 353 attention distribution based on input  
 354 features and their contextual impor-  
 355 tance, providing an adaptive refine-  
 356 ment of temporal relationships during  
 357 training—a capability absent in static  
 358 GCNs. This combination of stream-  
 359 lined design and enhanced aggrega-  
 360 tion yields an efficient, interpretable  
 361 model that bridges sequence and graph paradigms for superior  
 time series modeling.

362 **Remark 2.** Although Equation 10 presents a simplified formulation without multi-head attention,  
 363 residual connections, or layer normalization for clarity, Fighter is designed as a modular plug-in that  
 364 can seamlessly integrate these components into any attention-based framework. In our experimental  
 365 implementation, we adopt these standard enhancements to stabilize training and improve performance,  
 366 as is common in Transformer architectures. The Pseudocode 1 shows Fighter’s core forward pass,  
 367 with optional lines for incorporating multi-head attention (via parallel heads and concatenation),  
 368 residuals, and layer normalization.

371 5 EXPERIMENTS AND RESULTS

373 In our experiments, we evaluate the performance of Fighter on time series forecasting tasks to validate  
 374 its effectiveness in capturing temporal dependencies. The overall experimental results are presented  
 375 in Table 1, which clearly highlights the superiority of Fighter over standard Transformer-based  
 376 baselines including Autoformer (Wu et al., 2021), Informer(Zhou et al., 2021), Reformer (Kitaev  
 377 et al., 2020), and vanilla Transformer (Vaswani et al., 2017). Furthermore, Fighter exhibits versatility  
 in text sequence modeling, achieving competitive results with a streamlined architecture that enhances

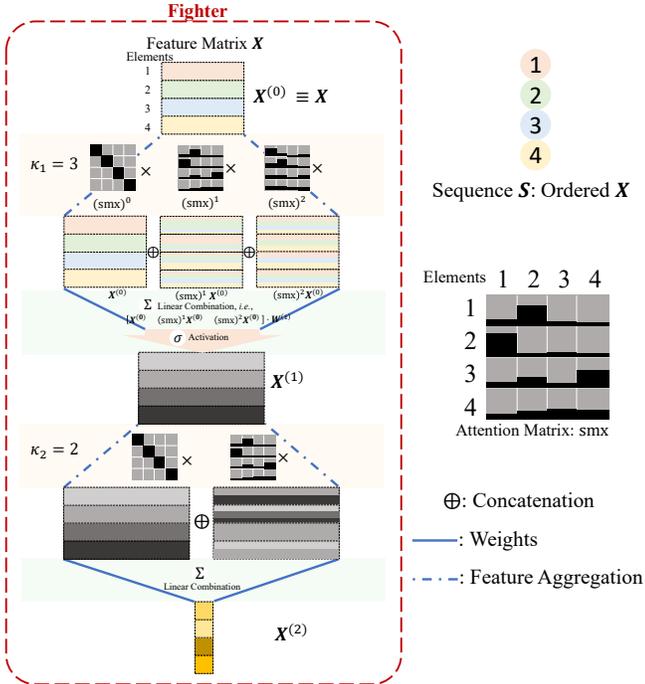


Figure 2: Visual illustration of the Fighter architecture, showcasing its streamlined design with multi-hop graph aggregation and dynamic adjacency matrix for efficient time series modeling.

Table 1: Time series forecasting (MSE/MAE) results on benchmark datasets and text sequence classification accuracy (Acc) on the AG News dataset. ‘O’ denotes the prediction length. Lower MSE/MAE indicates better forecasting performance, while higher accuracy reflects superior classification. The best results are highlighted in bold.

Dataset	O	Autoformer		Informer		Reformer		Transformer		Fighter ( $\kappa = 3$ )	
		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
Electricity	96	0.474±0.078	0.485±0.041	0.544±0.006	0.546±0.003	0.374±0.004	0.438±0.004	0.388±0.002	0.442±0.002	<b>0.339±0.007</b>	<b>0.417±0.004</b>
	192	0.581±0.157	0.533±0.074	0.652±0.062	0.616±0.032	0.387±0.010	0.447±0.007	0.384±0.019	<b>0.441±0.013</b>	<b>0.376±0.034</b>	0.448±0.027
	336	0.463±0.078	0.484±0.038	0.686±0.015	0.632±0.005	0.466±0.013	0.502±0.009	0.441±0.053	0.487±0.042	<b>0.408±0.015</b>	<b>0.477±0.012</b>
	720	0.559±0.194	0.539±0.105	0.966±0.007	0.801±0.004	0.490±0.014	0.519±0.008	0.442±0.045	0.484±0.033	<b>0.380±0.010</b>	<b>0.451±0.007</b>
Weather	96	0.392±0.023	0.444±0.013	0.911±0.095	0.709±0.039	1.076±0.129	0.813±0.054	0.564±0.005	0.545±0.004	<b>0.270±0.008</b>	<b>0.344±0.007</b>
	192	0.348±0.023	0.415±0.020	0.898±0.036	0.708±0.017	1.080±0.093	0.806±0.041	0.620±0.013	0.569±0.010	<b>0.256±0.012</b>	<b>0.330±0.009</b>
	336	0.384±0.037	0.430±0.028	0.930±0.020	0.729±0.010	1.038±0.019	0.788±0.002	0.669±0.102	0.585±0.050	<b>0.296±0.008</b>	<b>0.358±0.007</b>
	720	0.628±0.040	0.587±0.027	1.137±0.025	0.831±0.014	0.693±0.137	0.614±0.073	0.957±0.078	0.724±0.035	<b>0.342±0.004</b>	<b>0.383±0.004</b>
ETTh1	96	1.117±0.043	0.846±0.010	1.135±0.049	0.799±0.016	1.139±0.016	0.791±0.009	0.919±0.040	0.754±0.022	<b>0.732±0.076</b>	<b>0.651±0.050</b>
	192	1.148±0.085	0.839±0.026	1.108±0.060	0.793±0.028	1.251±0.035	0.840±0.014	0.949±0.016	0.762±0.007	<b>0.727±0.300</b>	<b>0.631±0.164</b>
	336	1.156±0.031	0.867±0.006	1.246±0.033	0.883±0.019	1.383±0.016	0.884±0.004	0.971±0.058	0.778±0.031	<b>0.912±0.242</b>	<b>0.759±0.135</b>
	720	1.108±0.044	0.846±0.019	1.267±0.019	0.873±0.010	1.425±0.009	0.894±0.003	1.049±0.007	0.810±0.009	<b>0.897±0.163</b>	<b>0.763±0.096</b>
ETTM1	96	0.903±0.004	0.736±0.007	0.916±0.044	0.718±0.027	1.034±0.015	0.744±0.009	0.572±0.046	0.542±0.024	<b>0.500±0.032</b>	<b>0.499±0.021</b>
	192	0.877±0.078	0.728±0.046	0.955±0.029	0.727±0.014	1.196±0.017	0.807±0.005	0.690±0.046	0.609±0.018	<b>0.517±0.036</b>	<b>0.514±0.034</b>
	336	0.999±0.014	0.788±0.031	1.046±0.068	0.771±0.031	1.391±0.018	0.885±0.006	0.807±0.014	0.675±0.010	<b>0.603±0.066</b>	<b>0.572±0.041</b>
	720	0.998±0.105	0.787±0.051	1.090±0.022	0.795±0.005	1.538±0.001	0.926±0.001	1.002±0.030	0.751±0.003	<b>0.701±0.026</b>	<b>0.627±0.029</b>
		Acc		Acc		Acc		Acc		Acc	
AG News	-	0.876		0.831		0.820		0.864		<b>0.898</b>	

efficiency compared to traditional Transformer-based baselines. Detailed settings and additional discussion are given in Appendix B.

We assess Fighter’s performance using widely recognized datasets for time series forecasting and text classification to ensure a comprehensive evaluation. These datasets are:

- Electricity<sup>3</sup>: containing the hourly electricity consumption of 321 customers from 2012 to 2014;
- Weather<sup>4</sup>: containing 21 meteorological indicators, recorded every 10 minutes for the 2020 year;
- ETT(Zhou et al., 2021): containing 7 electricity transformer factors, ETTh1 is recorded every hour, and ETTm1 is recorded every 15 minutes from 2016 to 2018;
- AG News(Zhang et al., 2015): containing 120,000 samples, with 30,000 instances per class, consisting of news articles from four categories: World, Sports, Business, and SCI / Technology.

Fighter extends the vanilla Transformer by incorporating higher-order compositions of the attention distribution matrix, controlled by the parameter  $\kappa$ , which governs the power of attention distributions. This mechanism enables **multi-hop graph aggregation**, explicitly modeling temporal dependencies across multiple scales as dynamic adjacency matrix, where dependencies are represented as graph edges. For the benchmark time series datasets, Fighter consistently outperforms most baselines, including Autoformer (Wu et al., 2021), Informer (Zhou et al., 2021), Reformer (Kitaev et al., 2020), and the vanilla Transformer (Vaswani et al., 2017) across different prediction lengths varying from 96 to 720. On the Electricity dataset, Fighter achieves the lowest MSE across all prediction horizons (96, 192, 336, and 720). Even at the longest horizon of 720, it records an MSE of 0.380, noticeably outperforming Transformer (0.442) and Reformer (0.490). At the horizon of 336, Fighter maintains a clear edge with an MSE of 0.408 compared with 0.441 (Transformer) and 0.466 (Reformer), while its MAE follows a similarly favorable trend, reflecting stable and reliable forecasting capability.

On the Weather dataset, Fighter exhibits particularly impressive results. For a prediction length of 96, it obtains an MSE/MAE of 0.270/0.344, which corresponds to a 74.9% lower MSE and 57.7% lower MAE than Reformer (1.076/0.813), and a reduction of approximately 31% and 23% relative to Autoformer (0.392/0.444). This advantageous gap persists across longer horizons; at  $O = 720$ , Fighter still achieves a competitive 0.342/0.383, remaining comfortably ahead of the other baselines and demonstrating remarkable robustness over extended sequences.

Fighter continues to perform favorably on the ETTh1 and ETTm1 datasets. At the longest tested horizon of 720 on ETTh1, it attains an MSE/MAE of 0.897/0.763, representing improvements of 37.1% and 14.6% over Reformer and 29.2% and 12.6% over Informer. The gains are even more pronounced on ETTm1 under the same setting, with reductions of 54.4% and 32.3% relative to Reformer, and 35.7% and 21.1% relative to Informer. Such sustained advantages as the prediction

<sup>3</sup><https://archive.ics.uci.edu/ml/datasets/ElectricityLoadDiagrams20112014>

<sup>4</sup><https://www.bgc-jena.mpg.de/wetter>

432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442

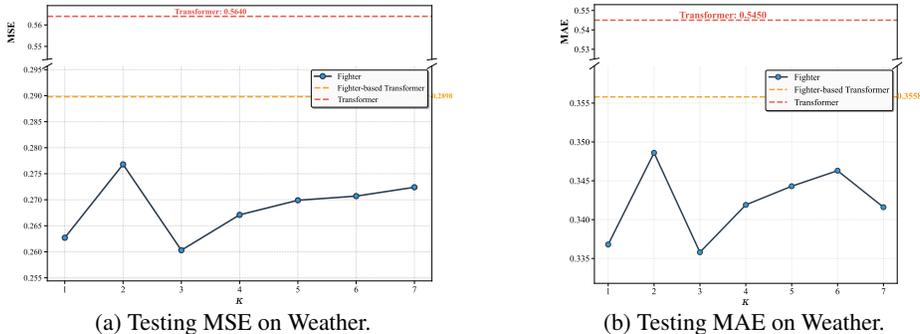


Figure 3: Time series forecasting performance of Fighter, Fighter-based Transformer with  $\kappa \in [1, 7]$  compared to the vanilla Transformer on the Weather dataset. Fighter-based Transformer denotes implementing the Transformer using the Fighter framework. Lower MSE/MAE indicate superior performance.

443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459

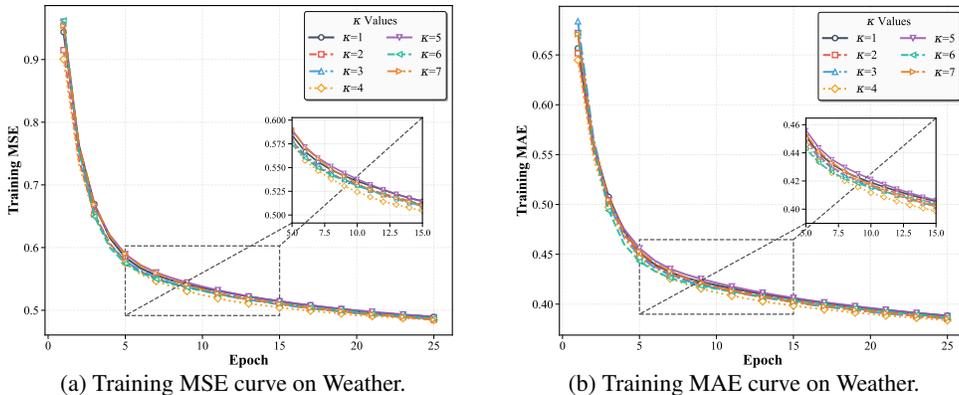


Figure 4: The training MSE/MAE performance of Fighter, with  $\kappa \in [1, 7]$  on the Weather dataset. The MSE also denotes the training loss. Lower MSE/MAE indicate superior performance.

length increases highlight Fighter’s effective modeling of long-range dependencies via its *multi-hop aggregation* mechanism.

Beyond time series forecasting, Fighter also shows promising versatility on the AG News text classification task, achieving 89.8% accuracy, which is higher than Transformer (86.4%), Autoformer (87.6%), and the remaining baselines, further illustrating its ability to handle diverse, high-dimensional sequential data effectively.

470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

Fighter’s multi-hop graph aggregation ensures better performance across varying prediction lengths, preventing rapid error accumulation. The interpretable nature of Fighter’s attention matrix, viewed as dynamic adjacency matrix, explains its superior performance: multi-hop paths allow long-range information to propagate naturally, reducing forecasting errors and enabling robust modeling of multivariate temporal relationships. To further explore the influence of the parameter  $\kappa$  on Fighter’s performance, we analyze its testing MSE/MAE on the Weather dataset, as presented in Figure 3. Overall, Fighter’s MSE and MAE surpass that of Fighter-based Transformer (implementing the Transformer using the Fighter framework) and the vanilla Transformer, meanwhile, Fighter and Fighter-based Transformer demonstrate better results across different  $\kappa$  than the vanilla Transformer. In detail, when  $\kappa$  is 3, Fighter obtains the best MSE/MAE results, which indicates Fighter’s multi-hop graph aggregation effectively captures long-range attention dependencies, leading to improved feature aggregation and significantly lower errors. When  $\kappa$  is 1 or 2, the dynamic attention distribution matrix captures only short-range dependencies, which is insufficient for modeling the Weather dataset’s complex temporal patterns. when  $\kappa$  is over 3, The MSE result remains a performance advantage over the two other baselines, indicating the sufficient capture of long-range dependencies. Beyond this point, increasing  $\kappa$  results in an potential over-smoothing effect (Keriven, 2022; Rusch et al., 2023), slightly degrading performance due to excessive aggregation. A similar pattern is observed in the testing MAE, confirming the critical role of  $\kappa$  in balancing dependency modeling. We also present the

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

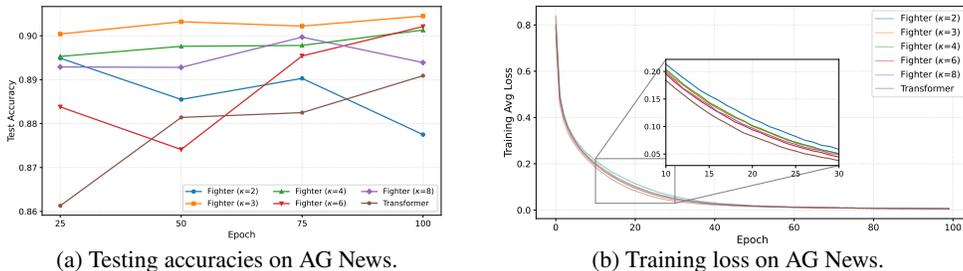


Figure 5: Text classification performance of Fighter with  $\kappa \in \{2, 3, 4, 6, 8\}$  versus the standard Transformer on the AG News dataset, evaluated across epochs 25 to 100. Higher accuracy reflects better performance.

training MSE/MAE performance comparison of Fighter across varying  $\kappa$  in Figure 4. Additionally, we examine  $\kappa$ 's effect for text sequence modeling on the AG News dataset, as shown in Figure 5.

To further validate the effectiveness of the proposed multi-hop graph aggregation mechanism, we conduct an ablation study by incorporating it ( $\kappa = 3$  applied on the attention distribution matrix) into the encoder-only iTransformer model. As shown in Table 2, integrating this mechanism consistently improves the forecasting performance of iTransformer on the ETTh1 and ETTm1 datasets across short and long horizons, with clear reductions in both MSE and MAE. These gains confirm that the multi-hop aggregation serves as a general and effective enhancement for capturing long-range dependencies, even when applied to strong state-of-the-art baselines.

Table 2: Ablation study on the multi-hop graph aggregation ( $\kappa = 3$ ). Forecasting performance (MSE/MAE) of the encoder-only iTransformer versus iTransformer + multi-hop graph aggregation on the ETTh1 and ETTm1 datasets at prediction horizons of 96 and 720. Best results are highlighted in bold.

Model	O	ETTh1		ETTM1	
		MSE	MAE	MSE	MAE
iTransformer	96	0.413	0.423	0.367	0.388
	720	0.547	0.516	0.497	0.460
iTransformer + Multi-hop ( $\kappa = 3$ )	96	<b>0.404</b>	<b>0.417</b>	<b>0.340</b>	<b>0.374</b>
	720	<b>0.522</b>	<b>0.505</b>	<b>0.488</b>	<b>0.454</b>

## 6 CONCLUDING REMARKS AND FUTURE WORK

This work establishes a novel equivalence between Transformer encoders and Graph Convolutional Networks (GCNs), showing that the attention distribution matrix serves as a dynamic adjacency matrix, with backward pass updates to the value matrix and FFN resembling that to GCN parameters. this unified reinterpretation demystifies Transformers in time series modeling, bridging sequence and graph paradigms. Our proposed **Fighter** architecture leverages this insight by removing redundant projections and integrating multi-hop graph aggregation, offering interpretable temporal dependency representations as graph edges. Fighter achieves competitive performance on forecasting benchmarks while enhancing mechanistic interpretability.

Future work would be extending this graph-convolutional reinterpretation to other Transformer variants such as Conformers (Gulati et al., 2020) in speech processing and Diffusion Transformers (Peebles & Xie, 2023) in vision and generative modeling. Such extensions would not only validate the generality of the sequence-graph equivalence across modalities, but also provide a principled framework to analyze and simplify diverse architectures, offering clearer interpretability and guiding the design of more efficient, domain-adaptive models.

## REPRODUCIBILITY STATEMENT

We have taken substantial measures to ensure the reproducibility of our work. The notations, theoretical derivations, and algorithms are thoroughly detailed in Appendix A. Appendix B provides

a complete description of the experimental setup, including training configurations, hyperparameter choices, and algorithmic details.

## STATEMENT ON THE USE OF LARGE LANGUAGE MODELS

We employed a large language model to enhance the manuscript’s language, such as improving grammar and phrasing. All research ideas, methods, experiments, analyses, figures/tables, and conclusions were solely developed by the authors, who take full responsibility for the content.

## REFERENCES

- Sami Abu-El-Haija, Bryan Perozzi, Amol Kapoor, Nazanin Alipourfard, Kristina Lerman, Hrayr Harutyunyan, Greg Ver Steeg, and Aram Galstyan. Mixhop: Higher-order graph convolutional architectures via sparsified neighborhood mixing. In *international conference on machine learning*, pp. 21–29. PMLR, 2019. 3
- Lei Bai, Lina Yao, Can Li, Xianzhi Wang, and Can Wang. Adaptive graph convolutional recurrent network for traffic forecasting. In *NeurIPS*, 2020. 2
- Defu Cao, Yujing Wang, Juanyong Duan, Ce Zhang, Xia Zhu, Congrui Huang, Yunhai Tong, Bixiong Xu, Jing Bai, Jie Tong, et al. Spectral temporal graph neural network for multivariate time-series forecasting. In *NeurIPS*, 2020. 2
- Ines Chami, Sami Abu-El-Haija, Bryan Perozzi, Christopher Ré, and Kevin Murphy. Machine learning on graphs: A model and comprehensive taxonomy. *Journal of Machine Learning Research*, 23 (89):1–64, 2022. 1, 5
- RG Cirstea, C Guo, B Yang, T Kieu, X Dong, and S Pan. Triformer: Triangular, variable-specific attentions for long sequence multivariate time series forecasting. In *The Thirty-First International Joint Conference on Artificial Intelligence (IJCAI-22)*. International Joint Conference on Artificial Intelligence (IJCAI), 2022. 2
- Izrail Moiseevitch Gelfand, Richard A Silverman, et al. *Calculus of variations*. Courier Corporation, 2000. 5
- Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al. Conformer: Convolution-augmented transformer for speech recognition. *arXiv preprint arXiv:2005.08100*, 2020. 10
- William L. Hamilton, Rex Ying, and Jure Leskovec. Representation learning on graphs: Methods and applications. *IEEE Data Engineering Bulletin*, 40:52–74, 2017. URL <https://api.semanticscholar.org/CorpusID:3215337>. 1
- Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012. 22
- Chaitanya K Joshi. Transformers are graph neural networks. *arXiv preprint arXiv:2506.22084*, 2025. 3, 4
- Nicolas Keriven. Not too little, not too much: a theoretical analysis of graph (over) smoothing. *Advances in Neural Information Processing Systems*, 35:2268–2281, 2022. 9
- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *ICLR*, 2017. 1, 2
- Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *International Conference on Learning Representations*, 2020. 7, 8
- Sanjukta Krishnagopal and Luana Ruiz. Graph neural tangent kernel: Convergence on large graphs. In *ICML*, 2023. 3

- 594 Yanhong Li and David C Anastasiu. Pfformer: A position-free transformer variant for extreme-  
595 adaptive multivariate time series forecasting. In *Pacific-Asia Conference on Knowledge Discovery*  
596 *and Data Mining*, pp. 185–197. Springer, 2025. 2
- 597
- 598 Yingcong Li, Yixiao Huang, Muhammed E Ildiz, Ankit Singh Rawat, and Samet Oymak. Mechanics  
599 of next token prediction with self-attention. In *International Conference on Artificial Intelligence*  
600 *and Statistics*, pp. 685–693. PMLR, 2024. 3, 4, 5
- 601 Shizhan Liu, Hang Yu, Cong Liao, Jianguo Li, Weiyao Lin, Alex X Liu, and Schahram Dust-  
602 dar. Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and  
603 forecasting. In *International Conference on Learning Representations*, 2022. 2
- 604
- 605 Ricardo P Masini, Marcelo C Medeiros, and Eduardo F Mendes. Machine learning advances for time  
606 series forecasting. *Journal of economic surveys*, 37(1):76–111, 2023. 1
- 607
- 608 Charles A Micchelli and Massimiliano Pontil. On learning vector-valued functions. *Neural computa-*  
609 *tion*, 17(1):177–204, 2005. 20
- 610 Mohammad Amin Morid, Olivia R Liu Sheng, and Joseph Dunbar. Time series prediction using deep  
611 learning methods in healthcare. *ACM Transactions on Management Information Systems*, 14(1):  
612 1–29, 2023. 1
- 613
- 614 Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64  
615 words: Long-term forecasting with transformers. In *The Eleventh International Conference on*  
616 *Learning Representations*, 2023. 2
- 617 Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier  
618 Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn:  
619 Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830, 2011. 22
- 620
- 621 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*  
622 *the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023. 10
- 623
- 624 Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint*  
625 *arXiv:1609.04747*, 2016. 5
- 626
- 627 Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use  
628 interpretable models instead. *Nature machine intelligence*, 1(5):206–215, 2019. 1
- 629
- 630 T Konstantin Rusch, Michael M Bronstein, and Siddhartha Mishra. A survey on oversmoothing in  
631 graph neural networks. *SAM Research Report*, 2023, 2023. 9
- 632
- 633 Chao Shang, Jie Chen, and Jinbo Bi. Discrete graph structure learning for forecasting multiple time  
634 series. In *ICLR*, 2021. 2
- 635
- 636 Tian Sun, Yuqi Chen, and Weiwei Sun. Penguin: Enhancing transformer with periodic-nested group  
637 attention for long-term time series forecasting. *arXiv preprint arXiv:2508.13773*, 2025. 2
- 638
- 639 Ruey S Tsay. *Analysis of financial time series*. John wiley & sons, 2005. 1
- 640
- 641 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
642 Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 1, 2, 3, 7, 8
- 643
- 644 Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua  
645 Bengio. Graph attention networks. In *ICLR*, 2018. 4
- 646
- 647 Eleni I Vlahogianni and Matthew G Karlaftis. Testing and comparing neural network and statistical  
approaches for predicting transportation time series. *Transportation research record*, 2399(1):  
9–22, 2013. 1
- Xue Wang, Jun-Jie Ma, Sheng Wang, and Dao-Wei Bi. Time series forecasting for energy-efficient  
organization of wireless sensor networks. *Sensors*, 7(9):1766–1792, 2007. 1

- 648 Qingsong Wen, Tian Zhou, Chaoli Zhang, Weiqi Chen, Ziqing Ma, Junchi Yan, and Liang Sun.  
649 Transformers in time series: a survey. In *Proceedings of the Thirty-Second International Joint*  
650 *Conference on Artificial Intelligence*, pp. 6778–6786, 2023. 1
- 651  
652 Kristoffer Wickstrøm, Karl Øyvind Mikalsen, Michael Kampffmeyer, Arthur Revhaug, and Robert  
653 Jenssen. Uncertainty-aware deep ensembles for reliable and explainable predictions of clinical  
654 time series. *IEEE Journal of Biomedical and Health Informatics*, 25(7):2435–2444, 2020. 1, 3
- 655 Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Simplify-  
656 ing graph convolutional networks. In *ICML*, 2019. 1, 2
- 657  
658 Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers  
659 with auto-correlation for long-term series forecasting. *Advances in neural information processing*  
660 *systems*, 34:22419–22430, 2021. 2, 7, 8
- 661 Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, Xiaojun Chang, and Chengqi Zhang. Connecting  
662 the dots: Multivariate time series forecasting with graph neural networks. In *KDD*, 2020. 2
- 663  
664 Kun Yi, Qi Zhang, Wei Fan, Hui He, Liang Hu, Pengyang Wang, Ning An, Longbing Cao, and  
665 Zhendong Niu. Fouriergnn: Rethinking multivariate time series forecasting from a pure graph  
666 perspective. In *NeurIPS*, 2023. 3
- 667 A Yarkın Yıldız, Emirhan Koç, and Aykut Koç. Multivariate time series imputation with transformers.  
668 *IEEE Signal Processing Letters*, 29:2517–2521, 2022. 3
- 669  
670 Bing Yu, Haoteng Yin, and Zhanxing Zhu. Spatio-temporal graph convolutional networks: A deep  
671 learning framework for traffic forecasting. In *Proceedings of the Twenty-Seventh International*  
672 *Joint Conference on Artificial Intelligence*, pp. 3634–3640. International Joint Conferences on  
673 Artificial Intelligence Organization, 2018. 2
- 674 Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series  
675 forecasting? In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp.  
676 11121–11128, 2023. 1, 2
- 677  
678 George Zerveas, Srideepika Jayaraman, Dhaval Patel, Anuradha Bhamidipaty, and Carsten Eickhoff.  
679 A transformer-based framework for multivariate time series representation learning. In *Proceedings*  
680 *of the 27th ACM SIGKDD conference on knowledge discovery & data mining*, pp. 2114–2124,  
681 2021. 3
- 682  
683 Chen Zhang, Xiaofeng Cao, Weiyang Liu, Ivor Tsang, and James Kwok. Nonparametric teaching for  
684 multiple learners. In *NeurIPS*, 2023. 20
- 685  
686 Chen Zhang, Weixin Bu, Zeyi Ren, Zhengwu Liu, Yik-Chung Wu, and Ngai Wong. Nonparametric  
687 teaching for graph property learners. In *ICML*, 2025. 1, 3
- 688  
689 Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text  
690 classification. *Advances in neural information processing systems*, 28, 2015. 8
- 691  
692 Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang.  
693 Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings*  
694 *of the AAAI conference on artificial intelligence*, volume 35, pp. 11106–11115, 2021. 2, 7, 8
- 695  
696 Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer: Frequency  
697 enhanced decomposed transformer for long-term series forecasting. In *International conference on*  
698 *machine learning*, pp. 27268–27286. PMLR, 2022. 2
- 699  
700  
701

# Appendix

## A ADDITIONAL DISCUSSIONS

### A.1 NOTATION OVERVIEW

Notation	Description
$\mathcal{G} \in \mathcal{G}$	Graph with $n$ nodes, node features $\mathbf{X}_{n \times d}$ , and adjacency matrix $\mathbf{A}_{n \times n}$
$\mathbf{X}_{n \times d}$	$n \times d$ node feature matrix; rows are node features $\mathbf{x}_i \in \mathbb{R}^d, i \in \mathbb{N}_n$
$\mathbf{A}_{n \times n}$	$n \times n$ adjacency matrix encoding edges
$\mathbb{N}_n := \{1, \dots, n\}$	Set of node indices
$\mathbb{N}_S := \{1, \dots, S\}$	Set of sequence indices
$\mathcal{S}_{1:S} \in \mathcal{S}$	Sequence of length $S$ with features $\mathbf{x}_s \in \mathbb{R}^d, s \in \mathbb{N}_S$
$\mathbf{X}_{S \times d}$	$S \times d$ sequence feature matrix
$\mathbf{x} = [x_1, \dots, x_d]$	Feature vector as row vector (transpose denotes column form)
$\mathbf{X}_{(i,:)}$	$i$ -th row of feature matrix $\mathbf{X}$ (i.e., $\mathbf{e}_i^\top \mathbf{X}$ )
$\mathbf{X}_{(:,j)}$	$j$ -th column of feature matrix $\mathbf{X}$ (i.e., $\mathbf{X} \mathbf{e}_j$ )
$\mathbf{e}_i$	Standard basis vector with 1 in $i$ -th position, 0 elsewhere
$\mathbf{I}$	Identity matrix
$\text{blkdiag}(a_1, \dots, a_m)$	Diagonal matrix with entries $a_1, \dots, a_m$
$\text{blkdiag}(a; m)$	Diagonal matrix with $m$ identical entries $a$ on the diagonal

Table 3: Summary of Key Notations.

### A.2 REFORMULATION OF A TRANSFORMER ENCODER

Omitting bias terms, Equation 3 for a Transformer encoder can be expressed as

$$\mathbf{X}_{\text{Att}}^{(\ell)} = \text{smx} \left( d^{-1/2} \mathbf{X}_{\text{FFN}}^{(\ell-1)} \mathbf{W}_Q^{(\ell)} \left( \mathbf{X}_{\text{FFN}}^{(\ell-1)} \mathbf{W}_K^{(\ell)} \right)^\top \right) \mathbf{X}_{\text{FFN}}^{(\ell-1)} \mathbf{W}_V^{(\ell)} \quad (11)$$

$$\mathbf{X}_{\text{FFN}}^{(\ell)} = \sigma \left( \mathbf{X}_{\text{Att}}^{(\ell)} \mathbf{W}_{\text{FFN}}^{(\ell,1)} \right) \mathbf{W}_{\text{FFN}}^{(\ell,2)} \quad (12)$$

$$\mathbf{X}_{\text{Att}}^{(L)} = \mathbf{X}_{\text{FFN}}^{(L-1)} \mathbf{W}^{(L)}, \quad (13)$$

Substituting Equation 11 into 12 yields:

$$\mathbf{X}_{\text{FFN}}^{(\ell)} = \sigma \left( \text{smx} \left( d^{-1/2} \mathbf{X}_{\text{FFN}}^{(\ell-1)} \mathbf{W}_Q^{(\ell)} \left( \mathbf{X}_{\text{FFN}}^{(\ell-1)} \mathbf{W}_K^{(\ell)} \right)^\top \right) \mathbf{X}_{\text{FFN}}^{(\ell-1)} \mathbf{W}_V^{(\ell)} \mathbf{W}_{\text{FFN}}^{(\ell,1)} \right) \mathbf{W}_{\text{FFN}}^{(\ell,2)} \quad (14)$$

Dropping the subscript for the hidden feature  $\mathbf{X}_{\text{FFN}}^{(\ell-1)}$  (denoted as  $\mathbf{X}^{(\ell-1)}$  for brevity) and defining the composite matrix  $\mathbf{W}_V^{(\ell,1)} := \mathbf{W}_V^{(\ell)} \mathbf{W}_{\text{FFN}}^{(\ell,1)}$ , this simplifies to:

$$\mathbf{X}^{(\ell)} = \sigma \left( \text{smx} \left( d^{-1/2} \mathbf{X}^{(\ell-1)} \mathbf{W}_Q^{(\ell)} \left( \mathbf{X}^{(\ell-1)} \mathbf{W}_K^{(\ell)} \right)^\top \right) \mathbf{X}^{(\ell-1)} \mathbf{W}_V^{(\ell,1)} \right) \mathbf{W}_{\text{FFN}}^{(\ell,2)}. \quad (15)$$

Thus, the Transformer encoder, derived from Equation 3, is:

$$\begin{aligned} \mathbf{X}^{(\ell)} &= \sigma \left( \text{smx} \left( d^{-1/2} \mathbf{X}^{(\ell-1)} \mathbf{W}_Q^{(\ell)} \left( \mathbf{X}^{(\ell-1)} \mathbf{W}_K^{(\ell)} \right)^\top \right) \mathbf{X}^{(\ell-1)} \mathbf{W}_V^{(\ell,1)} \right) \mathbf{W}_{\text{FFN}}^{(\ell,2)} \\ \mathbf{X}^{(L)} &= \mathbf{X}^{(L-1)} \mathbf{W}^{(L)}, \end{aligned} \quad (16)$$

### A.3 DERIVATION OF GRADIENTS FOR GRAPH CONVOLUTIONAL NETWORK

For a graph convolutional network (GCN), as defined by Equation 1, a two-layer architecture with multi-hop aggregation can be expressed as:

$$\mathbf{X}_{\text{GCN}}^{(2)} = \mathbf{A}^{[\kappa_2]} \text{blkdiag} \left( \sigma \left( \mathbf{A}^{[\kappa_1]} \text{blkdiag} \left( \mathbf{X}^{(0)}; \kappa_1 \right) \cdot \mathbf{W}^{(1)} \right); \kappa_2 \right) \cdot \mathbf{W}^{(2)}. \quad (17)$$

Using the chain rule, we compute the derivative of  $\mathbf{X}_{\text{GCN}}^{(2)}$  with respect to the second-layer weights  $\mathbf{W}^{(2)}$ , a vector of dimension  $\kappa_2 h_1$ , as follows:

$$\begin{aligned}
\frac{\partial \mathbf{X}_{\text{GCN}}^{(2)}}{\partial \mathbf{W}^{(2)}} &= \frac{\partial \mathbf{A}^{[\kappa_2]} \text{blkdiag}(\mathbf{X}^{(1)}; \kappa_2) \cdot \mathbf{W}^{(2)}}{\partial \mathbf{W}^{(2)}} \\
&= \mathbf{A}^{[\kappa_2]} \text{blkdiag}(\mathbf{X}^{(1)}; \kappa_2) \\
&= \underbrace{\mathbf{A}^{[\kappa_2]}}_{\text{size: } n \times \kappa_2 n} \underbrace{\text{blkdiag}(\sigma(\mathbf{A}^{[\kappa_1]} \text{blkdiag}(\mathbf{X}^{(0)}; \kappa_1) \mathbf{W}^{(1)}); \kappa_2)}_{\text{size: } \kappa_2 n \times \kappa_2 h_1}. \tag{18}
\end{aligned}$$

The derivative with respect to the first-layer weights is more involved. For  $i \in \mathbb{N}_{h_1}$

$$\begin{aligned}
\frac{\partial \mathbf{X}_{\text{GCN}}^{(2)}}{\partial \mathbf{W}_{(:,i)}^{(1)}} &= \frac{\partial \mathbf{A}^{[\kappa_2]} \text{blkdiag}(\mathbf{X}^{(1)} \mathbf{e}_i \mathbf{e}_i^\top; \kappa_2) \cdot \mathbf{W}^{(2)}}{\partial \mathbf{W}_{(:,i)}^{(1)}} \\
&= \frac{\partial \mathbf{A}^{[\kappa_2]} \underbrace{\text{blkdiag}(\mathbf{X}^{(1)} \mathbf{e}_i; \kappa_2)}_{\text{size: } \kappa_2 n \times \kappa_2} \cdot \underbrace{\text{blkdiag}(\mathbf{e}_i^\top; \kappa_2)}_{\text{size: } \kappa_2 \times \kappa_2 h_1} \mathbf{W}^{(2)}}{\partial \mathbf{W}_{(:,i)}^{(1)}} \\
&= \frac{\partial \mathbf{A}^{[\kappa_2]} \underbrace{\text{blkdiag}(\mathbf{X}^{(1)} \mathbf{e}_i; \kappa_2)}_{\text{size: } \kappa_2 n \times \kappa_2} \cdot \begin{pmatrix} \mathbf{W}_{(i-h_1+h_1)}^{(2)} \\ \dots \\ \mathbf{W}_{(i-h_1+\kappa_2 h_1)}^{(2)} \end{pmatrix}_{\text{size: } \kappa_2 \times 1}}{\partial \mathbf{W}_{(:,i)}^{(1)}} \\
&= \frac{\partial \mathbf{A}^{[\kappa_2]} \begin{pmatrix} \mathbf{X}^{(1)} \mathbf{e}_i \mathbf{W}_{(i-h_1+h_1)}^{(2)} \\ \dots \\ \mathbf{X}^{(1)} \mathbf{e}_i \mathbf{W}_{(i-h_1+\kappa_2 h_1)}^{(2)} \end{pmatrix}_{\text{size: } \kappa_2 n \times 1}}{\partial \mathbf{W}_{(:,i)}^{(1)}} \\
&= \mathbf{A}^{[\kappa_2]} \begin{pmatrix} \frac{\partial \mathbf{X}^{(1)} \mathbf{e}_i}{\partial \mathbf{W}_{(:,i)}^{(1)}} \mathbf{W}_{(i-h_1+h_1)}^{(2)} \\ \dots \\ \frac{\partial \mathbf{X}^{(1)} \mathbf{e}_i}{\partial \mathbf{W}_{(:,i)}^{(1)}} \mathbf{W}_{(i-h_1+\kappa_2 h_1)}^{(2)} \end{pmatrix} \\
&= \underbrace{\mathbf{A}^{[\kappa_2]}}_{\text{size: } n \times \kappa_2 n} \underbrace{\begin{pmatrix} \dot{\sigma} \cdot \mathbf{A}^{[\kappa_1]} \text{blkdiag}(\mathbf{X}^{(0)}; \kappa_1) \cdot \mathbf{W}_{(i-h_1+h_1)}^{(2)} \\ \dots \\ \dot{\sigma} \cdot \mathbf{A}^{[\kappa_1]} \text{blkdiag}(\mathbf{X}^{(0)}; \kappa_1) \cdot \mathbf{W}_{(i-h_1+\kappa_2 h_1)}^{(2)} \end{pmatrix}}_{\text{size: } \kappa_2 n \times \kappa_1 h_0}, \tag{19} \\
&\hspace{10em} \underbrace{\hspace{10em}}_{\text{size: } 1 \times \kappa_1 h_0}
\end{aligned}$$

where the derivative of the activation function is defined as  $\dot{\sigma} = \frac{\partial \sigma(\mathbf{A}^{[\kappa_1]} \text{blkdiag}(\mathbf{X}^{(0)}; \kappa_1) \mathbf{W}_{(:,i)}^{(1)})}{\partial \mathbf{A}^{[\kappa_1]} \text{blkdiag}(\mathbf{X}^{(0)}; \kappa_1) \mathbf{W}_{(:,i)}^{(1)}}$ .

When using the ReLU activation function, this simplifies to  $\dot{\sigma} \cdot \mathbf{A}^{[\kappa_1]} \text{blkdiag}(\mathbf{X}^{(0)}; \kappa_1) = \sigma(\mathbf{A}^{[\kappa_1]} \text{blkdiag}(\mathbf{X}^{(0)}; \kappa_1))$ .

Thus, the gradients with respect to  $\mathbf{W}^{(2)}$  and  $\mathbf{W}_{(:,i)}^{(1)}$  (for  $i \in \mathbb{N}_{h_1}$ ) are:

$$\begin{aligned} \frac{\partial \mathbf{X}_{\text{GCN}}^{(2)}}{\partial \mathbf{W}^{(2)}} &= \underbrace{\mathbf{A}^{[\kappa_2]} \text{blkdiag}(\sigma(\underbrace{\mathbf{A}^{[\kappa_1]} \text{blkdiag}(\mathbf{X}^{(0)}; \kappa_1) \mathbf{W}^{(1)}; \kappa_2)}_{1 \times \kappa_2 h_1}))}_{n \times \kappa_2 n}, \\ \frac{\partial \mathbf{X}_{\text{GCN}}^{(2)}}{\partial \mathbf{W}_{(:,i)}^{(1)}} &= \underbrace{\mathbf{A}^{[\kappa_2]} \left( \begin{array}{c} \dot{\sigma} \cdot \mathbf{A}^{[\kappa_1]} \text{blkdiag}(\mathbf{X}^{(0)}; \kappa_1) \cdot \mathbf{W}_{(i-h_1+h_1)}^{(2)} \\ \dots \\ \dot{\sigma} \cdot \mathbf{A}^{[\kappa_1]} \text{blkdiag}(\mathbf{X}^{(0)}; \kappa_1) \cdot \mathbf{W}_{(i-h_1+\kappa_2 h_1)}^{(2)} \end{array} \right)}_{1 \times \kappa_1 h_0}. \end{aligned} \quad (20)$$

When restricted to single-hop aggregation (*i.e.*, using only the  $\mathbf{A}^1$  term), the expressions simplify to:

$$\frac{\partial \mathbf{X}_{\text{GCN}}^{(2)}}{\partial \mathbf{W}^{(2)}} = \sigma(\mathbf{A} \mathbf{X}^{(0)} \mathbf{W}^{(1)}), \quad \frac{\partial \mathbf{X}_{\text{GCN}}^{(2)}}{\partial \mathbf{W}_{(:,i)}^{(1)}} = \dot{\sigma} \mathbf{A} \mathbf{X}^{(0)} \mathbf{W}_{(i,:)}^{(2)}. \quad (21)$$

Here,  $\mathbf{A}$  is a predefined adjacency matrix encoding the graph structure, and  $\mathbf{W}^{(\ell)}$  governs feature transformation.

#### A.4 DERIVATION OF GRADIENTS FOR TRANSFORMER ENCODER

For the Transformer encoder, we simplify Equation 4 by removing the unnecessary linear projection  $\mathbf{W}_{\text{FFN}}^{(\ell,2)}$  and fixing  $L = 2$ , resulting in the final layer output:

$$\mathbf{X}_{\text{TF}}^{(2)} = \sigma \left( \text{smx} \left( d^{-1/2} \mathbf{X}^{(0)} \mathbf{W}_Q^{(1)} (\mathbf{X}^{(0)} \mathbf{W}_K^{(1)})^\top \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \right) \cdot \mathbf{W}^{(2)}, \quad (22)$$

##### A.4.1 GRADIENTS W.R.T. FEATURE TRANSFORMATION PARAMETERS $\mathbf{W}^{(2)}$ AND $\mathbf{W}_{V(:,i)}^{(1,1)}$

The gradient with respect to  $\mathbf{W}^{(2)}$  is given by:

$$\frac{\partial \mathbf{X}_{\text{TF}}^{(2)}}{\partial \mathbf{W}^{(2)}} = \sigma \left( \text{smx} \left( d^{-1/2} \mathbf{X}^{(0)} \mathbf{W}_Q^{(1)} (\mathbf{X}^{(0)} \mathbf{W}_K^{(1)})^\top \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \right) \quad (23)$$

Similarly, the gradient with respect to  $\mathbf{W}_{V(:,i)}^{(1,1)}$  ( $i \in \mathbb{N}_{h_1}$ ) is

$$\begin{aligned} \frac{\partial \mathbf{X}_{\text{TF}}^{(2)}}{\partial \mathbf{W}_{V(:,i)}^{(1,1)}} &= \frac{\partial \sigma \left( \text{smx} \left( d^{-1/2} \mathbf{X}^{(0)} \mathbf{W}_Q^{(1)} (\mathbf{X}^{(0)} \mathbf{W}_K^{(1)})^\top \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \right) \cdot \mathbf{W}^{(2)}}{\partial \mathbf{W}_{V(:,i)}^{(1,1)}} \\ &= \dot{\sigma} \cdot \frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}^{(0)} \mathbf{W}_Q^{(1)} (\mathbf{X}^{(0)} \mathbf{W}_K^{(1)})^\top \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)}}{\partial \mathbf{W}_{V(:,i)}^{(1,1)}} \cdot \mathbf{W}^{(2)} \\ &= \dot{\sigma} \cdot \text{smx} \left( d^{-1/2} \mathbf{X}^{(0)} \mathbf{W}_Q^{(1)} (\mathbf{X}^{(0)} \mathbf{W}_K^{(1)})^\top \right) \mathbf{X}^{(0)} \mathbf{e}_i^\top \cdot \mathbf{W}^{(2)} \\ &= \dot{\sigma} \cdot \text{smx} \left( d^{-1/2} \mathbf{X}^{(0)} \mathbf{W}_Q^{(1)} (\mathbf{X}^{(0)} \mathbf{W}_K^{(1)})^\top \right) \mathbf{X}^{(0)} \cdot \mathbf{W}_{(i,:)}^{(2)}, \end{aligned} \quad (24)$$

where  $\dot{\sigma}$  denotes the derivative of the activation function.

## A.4.2 GRADIENTS WITH RESPECT TO QUERY AND KEY WEIGHT MATRICES

Computing the derivative of  $\mathbf{X}_{\text{TF}}^{(2)}$  with respect to the **query weight matrix** is more complex. For  $i \in \mathbb{N}_p$ ,

$$\begin{aligned}
\frac{\partial \mathbf{X}_{\text{TF}}^{(2)}}{\partial \mathbf{W}_{Q(:,i)}^{(1)}} &= \frac{\partial \sigma \left( \text{smx} \left( d^{-1/2} \mathbf{X}^{(0)} \mathbf{W}_Q^{(1)} \left( \mathbf{X}^{(0)} \mathbf{W}_K^{(1)} \right)^\top \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \right) \cdot \mathbf{W}^{(2)}}{\partial \mathbf{W}_{Q(:,i)}^{(1)}} \\
&= \dot{\sigma} \frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}^{(0)} \mathbf{W}_Q^{(1)} \left( \mathbf{X}^{(0)} \mathbf{W}_K^{(1)} \right)^\top \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)}}{\partial \mathbf{W}_{Q(:,i)}^{(1)}} \cdot \mathbf{W}^{(2)} \\
&= \dot{\sigma} \frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{e}_i \mathbf{e}_i^\top \left( \mathbf{X}^{(0)} \mathbf{W}_K^{(1)} \right)^\top \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)}}{\partial \mathbf{W}_{Q(:,i)}^{(1)}} \cdot \mathbf{W}^{(2)} \\
&= \begin{bmatrix} \dot{\sigma}_1 \frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}_{(1,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)}}{\partial \mathbf{W}_{Q(:,i)}^{(1)}} \cdot \mathbf{W}^{(2)} \\ \dots \\ \dot{\sigma}_S \frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}_{(S,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)}}{\partial \mathbf{W}_{Q(:,i)}^{(1)}} \cdot \mathbf{W}^{(2)} \end{bmatrix} \\
&= \begin{bmatrix} \dot{\sigma}_1 \frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}_{(1,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right)}{\partial \mathbf{W}_{Q(:,i)}^{(1)}} \cdot \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)} \\ \dots \\ \dot{\sigma}_S \frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}_{(S,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right)}{\partial \mathbf{W}_{Q(:,i)}^{(1)}} \cdot \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)} \end{bmatrix} \quad (25)
\end{aligned}$$

To unpack this further, we analyze it row by row. For the  $j$ -th row ( $j \in \mathbb{N}_S$ ) in Equation 25, the key

derivative term is  $\frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right)}{\partial \mathbf{W}_{Q(:,i)}^{(1)}}$

$$\begin{aligned}
&\frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right)}{\partial \mathbf{W}_{Q(:,i)}^{(1)}} \\
&= \frac{\partial d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \partial \text{smx} \left( d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right)}{\partial \mathbf{W}_{Q(:,i)}^{(1)} \partial d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top}} \\
&= \mathbf{X}_{(j,:)}^{(0)} \cdot \frac{1}{\sqrt{d}} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right)}{\partial d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top}} \\
&= \mathbf{X}_{(j,:)}^{(0)} \cdot \frac{1}{\sqrt{d}} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \frac{\partial \left( \frac{\exp \left( d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right)}{\mathbf{1}^\top \exp \left( d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right)} \right)}{\partial d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top}} \\
&= \mathbf{X}_{(j,:)}^{(0)} \cdot \left( \frac{1}{\sqrt{d}} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \text{blkdiag} \left( \text{smx} \left( \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} / \sqrt{d} \right) \right. \right. \\
&\quad \left. \left. - \frac{1}{\sqrt{d}} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \left( \text{smx} \left( \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} / \sqrt{d} \right) \right)^\top \text{smx} \left( \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} / \sqrt{d} \right) \right) \quad (26)
\end{aligned}$$

The right-hand side of Equation 26 simplifies to:

$$\begin{aligned}
& \underbrace{\mathbf{X}_{(j,:)}^{(0)}}_{\text{size: } 1 \times d} \cdot \left( \underbrace{d^{-1/2} \mathbf{W}_{K(:,i)}^{(1) \top}}_{1 \times 1} \underbrace{\mathbf{X}^{(0) \top}}_{1 \times d} \underbrace{\text{blkdiag} \left( \underbrace{\text{smx}(\mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)})}_{1 \times d} \underbrace{\mathbf{W}_{K(:,i)}^{(1) \top}}_{d \times 1} \underbrace{\mathbf{X}^{(0) \top}}_{1 \times d} \underbrace{/\sqrt{d}}_{d \times S} \right)}_{S \times S} \right) \\
& - \underbrace{d^{-1/2} \mathbf{W}_{K(:,i)}^{(1) \top}}_{1 \times 1} \underbrace{\mathbf{X}^{(0) \top}}_{1 \times d} \left( \underbrace{\text{smx}(\mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)})}_{1 \times d} \underbrace{\mathbf{W}_{K(:,i)}^{(1) \top}}_{d \times 1} \underbrace{\mathbf{X}^{(0) \top}}_{1 \times d} \underbrace{/\sqrt{d}}_{d \times S} \right) \underbrace{\text{smx}(\mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)})}_{1 \times d} \underbrace{\mathbf{W}_{K(:,i)}^{(1) \top}}_{d \times 1} \underbrace{\mathbf{X}^{(0) \top}}_{1 \times d} \underbrace{/\sqrt{d}}_{d \times S} \\
& = \underbrace{\mathbf{X}_{(j,:)}^{(0)}}_{1 \times d} \cdot \left( \underbrace{d^{-1/2} \mathcal{K}_{(:,i)}^{(1) \top}}_{1 \times 1} \underbrace{\text{blkdiag} \left( \underbrace{\text{smx}(\mathcal{Q}_{(j,i)}^{(1)} \mathcal{K}_{(:,i)}^{(1) \top})}_{1 \times S} \right)}_{S \times S} \right) \\
& - \underbrace{d^{-1/2} \mathcal{K}_{(:,i)}^{(1) \top}}_{1 \times 1} \underbrace{\left( \underbrace{\text{smx}(\mathcal{Q}_{(j,i)}^{(1)} \mathcal{K}_{(:,i)}^{(1) \top})}_{S \times 1} \right) \underbrace{\text{smx}(\mathcal{Q}_{(j,i)}^{(1)} \mathcal{K}_{(:,i)}^{(1) \top})}_{1 \times S}}_{1 \times S} \\
& = \underbrace{d^{-1/2} \mathbf{X}_{(j,:)}^{(0)}}_{1 \times d} \cdot \left( \underbrace{\mathcal{K}_{(:,i)}^{(1) \top}}_{1 \times S} \underbrace{\text{blkdiag} \left( \underbrace{\text{smx}(\xi_{(j,i;:,i)}^{(1)})}_{S \times S} \right)}_{S \times S} - \underbrace{\mathcal{K}_{(:,i)}^{(1) \top}}_{1 \times S} \underbrace{\left( \underbrace{\text{smx}(\xi_{(j,i;:,i)}^{(1)})}_{S \times S} \right) \underbrace{\text{smx}(\xi_{(j,i;:,i)}^{(1)})}_{S \times S}}_{1 \times S} \right), \tag{27}
\end{aligned}$$

where  $\mathcal{Q}^{(\ell)} := \mathbf{X}^{(\ell-1)} \mathbf{W}_Q^{(\ell)}$ ,  $\mathcal{K}^{(\ell)} := \mathbf{X}^{(\ell-1)} \mathbf{W}_K^{(\ell)}$ , and  $\xi_{(j,i;:,i)}^{(1)} := \mathcal{Q}_{(j,i)}^{(1)} \mathcal{K}_{(:,i)}^{(1) \top} / \sqrt{d}$ .

Integrating Equations 25, 26, and 27, and dropping the superscripts for  $\mathcal{Q}$  and  $\mathcal{K}$  for brevity, yields:

$$\begin{aligned}
& \frac{\partial \mathbf{X}_{\text{TF}}^{(2)}}{\partial \mathbf{W}_{Q(:,i)}^{(1)}} \\
& = \left[ \underbrace{\dot{\sigma}_1 / \sqrt{d} \mathbf{X}_{(1,:)}^{(0)}}_{1 \times d} \cdot \underbrace{\left( \underbrace{\mathcal{K}_{(:,i)}^{(1) \top}}_{1 \times S} \underbrace{\text{blkdiag} \left( \underbrace{\text{smx}(\xi_{(1,i;:,i)}^{(1)})}_{S \times S} \right)}_{S \times S} \mathbf{X}^{(0)} \underbrace{\mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)}}_{d \times h_1 \quad h_1 \times 1} - \underbrace{\mathcal{K}_{(:,i)}^{(1) \top}}_{1 \times S} \left( \underbrace{\text{smx}(\xi_{(1,i;:,i)}^{(1)})}_{S \times S} \right) \underbrace{\text{smx}(\xi_{(1,i;:,i)}^{(1)})}_{S \times d} \underbrace{\mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)}}_{h_1 \times 1} \right)}_{S \times d} \\
& \quad \underbrace{\dot{\sigma}_S / \sqrt{d} \mathbf{X}_{(S,:)}^{(0)}}_{1 \times d} \cdot \underbrace{\left( \underbrace{\mathcal{K}_{(:,i)}^{(1) \top}}_{1 \times S} \underbrace{\text{blkdiag} \left( \underbrace{\text{smx}(\xi_{(S,i;:,i)}^{(1)})}_{S \times S} \right)}_{S \times S} \mathbf{X}^{(0)} \underbrace{\mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)}}_{d \times h_1 \quad h_1 \times 1} - \underbrace{\mathcal{K}_{(:,i)}^{(1) \top}}_{1 \times S} \left( \underbrace{\text{smx}(\xi_{(S,i;:,i)}^{(1)})}_{S \times S} \right) \underbrace{\text{smx}(\xi_{(S,i;:,i)}^{(1)})}_{S \times d} \underbrace{\mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)}}_{h_1 \times 1} \right)}_{S \times d} \\
& = \left[ \dot{\sigma}_j / \sqrt{d} \mathbf{X}_{(j,:)}^{(0)} \cdot \left( \mathcal{K}_{(:,i)}^{(1) \top} \text{blkdiag} \left( \text{smx}(\xi_{(j,i;:,i)}^{(1)}) \right) - \mathcal{K}_{(:,i)}^{(1) \top} \left( \text{smx}(\xi_{(j,i;:,i)}^{(1)}) \right) \text{smx}(\xi_{(j,i;:,i)}^{(1)}) \right) \cdot \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)} \right]_{S \times d} \\
& = \text{blkdiag} \left( \dot{\sigma}_1 / \sqrt{d} \left( \mathcal{K}_{(:,i)}^{(1) \top} \text{blkdiag} \left( \text{smx}(\xi_{(1,i;:,i)}^{(1)}) \right) - \mathcal{K}_{(:,i)}^{(1) \top} \left( \text{smx}(\xi_{(1,i;:,i)}^{(1)}) \right) \text{smx}(\xi_{(1,i;:,i)}^{(1)}) \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)}, \right. \\
& \quad \left. \dots, \dot{\sigma}_S / \sqrt{d} \left( \mathcal{K}_{(:,i)}^{(1) \top} \text{blkdiag} \left( \text{smx}(\xi_{(S,i;:,i)}^{(1)}) \right) - \mathcal{K}_{(:,i)}^{(1) \top} \left( \text{smx}(\xi_{(S,i;:,i)}^{(1)}) \right) \text{smx}(\xi_{(S,i;:,i)}^{(1)}) \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)} \right) \mathbf{X}^{(0)}. \tag{28}
\end{aligned}$$

The derivative with respect to the **key weight matrix** follows a parallel approach. For  $i \in \mathbb{N}_p$ ,

$$\begin{aligned}
\frac{\partial \mathbf{X}_{\text{TF}}^{(2)}}{\partial \mathbf{W}_{K(:,i)}^{(1)}} &= \frac{\partial \sigma \left( \text{smx} \left( d^{-1/2} \mathbf{X}^{(0)} \mathbf{W}_Q^{(1)} (\mathbf{X}^{(0)} \mathbf{W}_K^{(1)})^\top \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \right) \cdot \mathbf{W}^{(2)}}{\partial \mathbf{W}_{K(:,i)}^{(1)}} \\
&= \dot{\sigma} \frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}^{(0)} \mathbf{W}_Q^{(1)} (\mathbf{X}^{(0)} \mathbf{W}_K^{(1)})^\top \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)}}{\partial \mathbf{W}_{K(:,i)}^{(1)}} \cdot \mathbf{W}^{(2)} \\
&= \dot{\sigma} \frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}^{(0)} \mathbf{W}_Q^{(1)} (\mathbf{X}^{(0)} \mathbf{W}_K^{(1)} \mathbf{e}_i \mathbf{e}_i^\top)^\top \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)}}{\partial \mathbf{W}_{K(:,i)}^{(1)}} \cdot \mathbf{W}^{(2)} \\
&= \begin{bmatrix} \dot{\sigma}_1 \frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}_{(1,:)}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)}}{\partial \mathbf{W}_{K(:,i)}^{(1)}} \cdot \mathbf{W}^{(2)} \\ \dots \\ \dot{\sigma}_S \frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}_{(S,:)}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)}}{\partial \mathbf{W}_{K(:,i)}^{(1)}} \cdot \mathbf{W}^{(2)} \end{bmatrix} \\
&= \begin{bmatrix} \dot{\sigma}_1 \frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}_{(1,:)}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right)}{\partial \mathbf{W}_{K(:,i)}^{(1)}} \cdot \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)} \\ \dots \\ \dot{\sigma}_S \frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}_{(S,:)}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right)}{\partial \mathbf{W}_{K(:,i)}^{(1)}} \cdot \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)} \end{bmatrix} \quad (29)
\end{aligned}$$

Again, examining row by row for the  $j$ -th row ( $j \in \mathbb{N}_S$ ) in Equation 29, the focal derivative is

$$\begin{aligned}
&\frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right)}{\partial \mathbf{W}_{K(:,i)}^{(1)}}: \\
&\frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right)}{\partial \mathbf{W}_{K(:,i)}^{(1)}} \\
&= \frac{\partial d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \partial \text{smx} \left( d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right)}{\partial \mathbf{W}_{K(:,i)}^{(1)} \partial d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top}} \\
&= \mathbf{X}_{(j,:)}^{(0)} \cdot \frac{1}{\sqrt{d}} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \frac{\partial \text{smx} \left( d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right)}{\partial d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top}} \\
&= \mathbf{X}_{(j,:)}^{(0)} \cdot \frac{1}{\sqrt{d}} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \frac{\partial \left( \frac{\exp \left( d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right)}{\mathbf{1}^\top \exp \left( d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \right)} \right)}{\partial d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top}} \\
&= \mathbf{X}_{(j,:)}^{(0)} \cdot \left( \frac{1}{\sqrt{d}} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \text{blkdiag} \left( \text{smx} \left( \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} / \sqrt{d} \right) \right. \right. \\
&\quad \left. \left. - \frac{1}{\sqrt{d}} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} \left( \text{smx} \left( \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} / \sqrt{d} \right) \right)^\top \text{smx} \left( \mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_Q^{(1)} \mathbf{W}_{K(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1)\top} \mathbf{X}^{(0)\top} / \sqrt{d} \right) \right) \quad (30)
\end{aligned}$$

Simplifying the right-hand side of Equation 30 gives:

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

$$\begin{aligned}
& \mathbf{X}_{(j,:)}^{(0)} \cdot \left( \underbrace{d^{-1/2} \mathbf{W}_{Q(:,i)}^{(1) \top}}_{1 \times 1} \underbrace{\mathbf{X}^{(0) \top}}_{d \times S} \text{blkdiag} \left( \underbrace{\text{smx}(\mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1) \top} \mathbf{X}^{(0) \top} / \sqrt{d})}_{S \times S} \right) \right. \\
& \left. - \underbrace{d^{-1/2} \mathbf{W}_{Q(:,i)}^{(1) \top}}_{1 \times 1} \underbrace{\mathbf{X}^{(0) \top}}_{1 \times d} \left( \underbrace{\text{smx}(\mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1) \top} \mathbf{X}^{(0) \top} / \sqrt{d})}_{S \times 1} \right) \underbrace{\text{smx}(\mathbf{X}_{(j,:)}^{(0)} \mathbf{W}_{Q(:,i)}^{(1)} \mathbf{W}_{K(:,i)}^{(1) \top} \mathbf{X}^{(0) \top} / \sqrt{d})}_{1 \times S} \right) \\
& = \mathbf{X}_{(j,:)}^{(0)} \cdot \left( \underbrace{d^{-1/2} \mathcal{Q}_{(:,i)}^{(1) \top}}_{1 \times 1} \underbrace{\mathbf{X}^{(0) \top}}_{1 \times S} \text{blkdiag} \left( \underbrace{\text{smx}(\mathcal{Q}_{(j,i)}^{(1)} \mathcal{K}_{(:,i)}^{(1) \top} / \sqrt{d})}_{S \times S} \right) \right. \\
& \left. - \underbrace{d^{-1/2} \mathcal{Q}_{(:,i)}^{(1) \top}}_{1 \times 1} \underbrace{\mathbf{X}^{(0) \top}}_{1 \times S} \left( \underbrace{\text{smx}(\mathcal{Q}_{(j,i)}^{(1)} \mathcal{K}_{(:,i)}^{(1) \top} / \sqrt{d})}_{S \times 1} \right) \underbrace{\text{smx}(\mathcal{Q}_{(j,i)}^{(1)} \mathcal{K}_{(:,i)}^{(1) \top} / \sqrt{d})}_{1 \times S} \right) \\
& = d^{-1/2} \mathbf{X}_{(j,:)}^{(0)} \cdot \left( \underbrace{\mathcal{Q}_{(:,i)}^{(1) \top}}_{1 \times 1} \text{blkdiag} \left( \underbrace{\text{smx}(\xi_{(j,i;:,i)}^{(1)})}_{S \times S} \right) - \underbrace{\mathcal{Q}_{(:,i)}^{(1) \top}}_{1 \times S} \left( \underbrace{\text{smx}(\xi_{(j,i;:,i)}^{(1)})}_{S \times S} \right) \underbrace{\text{smx}(\xi_{(j,i;:,i)}^{(1)})}_{S \times S} \right), \tag{31}
\end{aligned}$$

Combining Equations 29, 30, and 31 results in:

$$\begin{aligned}
& \frac{\partial \mathbf{X}_{\text{TF}}^{(2)}}{\partial \mathbf{W}_{K(:,i)}^{(1)}} \\
& = \left[ \begin{array}{c} \dot{\sigma}_1 / \sqrt{d} \mathbf{X}_{(1,:)}^{(0)} \cdot \left( \underbrace{\mathcal{Q}_{(:,i)}^{(1) \top}}_{1 \times 1} \text{blkdiag} \left( \underbrace{\text{smx}(\xi_{(1,i;:,i)}^{(1)})}_{S \times S} \right) \mathbf{X}^{(0)} \underbrace{\mathbf{W}_V^{(1,1)}}_{S \times d} \underbrace{\mathbf{W}^{(2)}}_{d \times h_1} \underbrace{\mathbf{W}^{(2)}}_{h_1 \times 1} - \underbrace{\mathcal{Q}_{(:,i)}^{(1) \top}}_{1 \times 1} \left( \underbrace{\text{smx}(\xi_{(1,i;:,i)}^{(1)})}_{S \times S} \right) \underbrace{\text{smx}(\xi_{(1,i;:,i)}^{(1)})}_{S \times d} \underbrace{\mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)}}_{d \times h_1} \underbrace{\mathbf{W}^{(2)}}_{h_1 \times 1} \right) \\ \dots \\ \dot{\sigma}_S / \sqrt{d} \mathbf{X}_{(S,:)}^{(0)} \cdot \left( \underbrace{\mathcal{Q}_{(:,i)}^{(1) \top}}_{1 \times 1} \text{blkdiag} \left( \underbrace{\text{smx}(\xi_{(S,i;:,i)}^{(1)})}_{S \times S} \right) \mathbf{X}^{(0)} \underbrace{\mathbf{W}_V^{(1,1)}}_{S \times d} \underbrace{\mathbf{W}^{(2)}}_{d \times h_1} \underbrace{\mathbf{W}^{(2)}}_{h_1 \times 1} - \underbrace{\mathcal{Q}_{(:,i)}^{(1) \top}}_{1 \times 1} \left( \underbrace{\text{smx}(\xi_{(S,i;:,i)}^{(1)})}_{S \times S} \right) \underbrace{\text{smx}(\xi_{(S,i;:,i)}^{(1)})}_{S \times d} \underbrace{\mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)}}_{d \times h_1} \underbrace{\mathbf{W}^{(2)}}_{h_1 \times 1} \right) \end{array} \right]_{S \times d} \\
& = \left[ \dot{\sigma}_j / \sqrt{d} \mathbf{X}_{(j,:)}^{(0)} \cdot \left( \mathcal{Q}_{(:,i)}^{(1) \top} \text{blkdiag} \left( \text{smx}(\xi_{(j,i;:,i)}^{(1)}) \right) - \mathcal{Q}_{(:,i)}^{(1) \top} \left( \text{smx}(\xi_{(j,i;:,i)}^{(1)}) \right) \text{smx}(\xi_{(j,i;:,i)}^{(1)}) \right) \cdot \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)} \right]_{S \times d} \\
& = \text{blkdiag} \left( \dot{\sigma}_1 / \sqrt{d} \left( \mathcal{Q}_{(:,i)}^{(1) \top} \text{blkdiag} \left( \text{smx}(\xi_{(1,i;:,i)}^{(1)}) \right) - \mathcal{Q}_{(:,i)}^{(1) \top} \left( \text{smx}(\xi_{(1,i;:,i)}^{(1)}) \right) \text{smx}(\xi_{(1,i;:,i)}^{(1)}) \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)}, \right. \\
& \quad \left. \dots, \dot{\sigma}_S / \sqrt{d} \left( \mathcal{Q}_{(:,i)}^{(1) \top} \text{blkdiag} \left( \text{smx}(\xi_{(S,i;:,i)}^{(1)}) \right) - \mathcal{Q}_{(:,i)}^{(1) \top} \left( \text{smx}(\xi_{(S,i;:,i)}^{(1)}) \right) \text{smx}(\xi_{(S,i;:,i)}^{(1)}) \right) \mathbf{X}^{(0)} \mathbf{W}_V^{(1,1)} \mathbf{W}^{(2)} \right) \mathbf{X}^{(0)}. \tag{32}
\end{aligned}$$

These derivations naturally extend to scenarios where each component of the output  $\mathbf{X}_{\text{TF}}^{(L)}$  is a vector, through a multi-dimensional framework (Micchelli & Pontil, 2005; Zhang et al., 2023). For multi-head attention, the process can be parallelized by replicating the derivations across the number of heads.

## A.5 FIGHTER FORWARD PASS (WITH OPTIONAL ENHANCEMENTS)

**Algorithm 1** Fighter Forward Pass (with Optional Enhancements)

---

**Require:** Input features  $\mathbf{X}_{\text{FT}}^{(0)} \in \mathbb{R}^{S \times d}$ , temporary variable  $\mathbf{A}$ , number of layers  $L$ , hop parameters  $\{\kappa_\ell\}_{\ell=1}^L$ , weights  $\{\mathbf{W}_Q^{(\ell)}, \mathbf{W}_K^{(\ell)}, \mathbf{W}^{(\ell)}\}_{\ell=1}^L$

**Ensure:** Output  $\mathbf{X}_{\text{FT}}^{(L)}$

- 1:  $\mathbf{X} \leftarrow \mathbf{X}_{\text{FT}}^{(0)}$  ▷ Initialize with input features
- 2: **for**  $\ell = 1$  to  $L - 1$  **do**
- 3:    *(Optional):*  $\tilde{\mathbf{X}} \leftarrow \text{LayerNorm}(\mathbf{X})$  ▷ Apply layer normalization
- 4:     $\mathbf{Q} \leftarrow \tilde{\mathbf{X}} \mathbf{W}_Q^{(\ell)}$  ▷ Queries; use  $\mathbf{X}$  if no norm
- 5:     $\mathbf{K} \leftarrow \tilde{\mathbf{X}} \mathbf{W}_K^{(\ell)}$  ▷ Keys; use  $\mathbf{X}$  if no norm
- 6:    *(Optional):* Compute  $\mathbf{Q}, \mathbf{K}$  for each head, concatenate/project outputs
- 7:     $\mathbf{A} \leftarrow \text{smx}(d^{-1/2} \mathbf{Q} \mathbf{K}^\top)$  ▷ Compute attention distribution matrices
- 8:     $\tilde{\mathbf{A}} \leftarrow \mathbf{A}^{[\kappa_\ell]}$  ▷ Multi-hop attention: raise to powers up to  $\kappa_\ell - 1$
- 9:     $\tilde{\tilde{\mathbf{X}}} \leftarrow \text{blkdiag}(\tilde{\mathbf{X}}; \kappa_\ell)$  ▷ Block-diagonal replication; use  $\mathbf{X}$  if no norm
- 10:     $\tilde{\mathbf{X}} \leftarrow \tilde{\tilde{\mathbf{A}}} \tilde{\tilde{\mathbf{X}}} \mathbf{W}^{(\ell)}$  ▷ Aggregate and project with streamlined weight
- 11:     $\mathbf{X} \leftarrow \sigma(\tilde{\mathbf{X}})$  ▷ Apply activation (e.g., ReLU)
- 12:    *(Optional):*  $\mathbf{X} \leftarrow \mathbf{X} + \tilde{\mathbf{X}}$  ▷ Add residual; use  $\mathbf{X}$  if no norm
- 13: **end for**
- 14: *(Optional):*  $\tilde{\mathbf{X}} \leftarrow \text{LayerNorm}(\mathbf{X})$  ▷ Final layer normalization
- 15:  $\mathbf{Q} \leftarrow \tilde{\mathbf{X}} \mathbf{W}_Q^{(L)}$  ▷ Final queries; use  $\mathbf{X}$  if no norm
- 16:  $\mathbf{K} \leftarrow \tilde{\mathbf{X}} \mathbf{W}_K^{(L)}$  ▷ Final keys; use  $\mathbf{X}$  if no norm
- 17: *(Optional):* Compute  $\mathbf{Q}, \mathbf{K}$  for each head, concatenate/project outputs
- 18:  $\mathbf{A} \leftarrow \text{smx}(d^{-1/2} \mathbf{Q} \mathbf{K}^\top)$  ▷ Final attention distribution matrices
- 19:  $\tilde{\mathbf{A}} \leftarrow \mathbf{A}^{[\kappa_L]}$  ▷ Final multi-hop attention
- 20:  $\tilde{\tilde{\mathbf{X}}} \leftarrow \text{blkdiag}(\tilde{\mathbf{X}}; \kappa_L)$  ▷ Final block-diagonal replication; use  $\mathbf{X}$  if no norm
- 21:  $\mathbf{X}_{\text{FT}}^{(L)} \leftarrow \tilde{\tilde{\mathbf{A}}}^{[\kappa_L]} \tilde{\tilde{\mathbf{X}}} \mathbf{W}^{(L)}$  ▷ Final output
- 22: *(Optional):*  $\mathbf{X}_{\text{FT}}^{(L)} \leftarrow \mathbf{X}_{\text{FT}}^{(L)} + \tilde{\mathbf{X}}$  ▷ Add final residual; use  $\mathbf{X}$  if no norm **return**  $\mathbf{X}_{\text{FT}}^{(L)}$

---

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

## B EXPERIMENT DETAILS

This section outlines the experiment details, covering the experimental setup, supplementary results, and a brief sensitivity analysis on benchmark datasets.

### B.1 EXPERIMENTAL SETUP

**Device Setup.** We mainly conduct experiments using NVIDIA Geforce RTX 3090 (24G).

**Dataset Splitting.** The train / val / test split configurations for the benchmark datasets are provided in Table 4. We adopt the Electricity, Weather and ETTh1, ETTm1 datasets for time series forecasting. All time series inputs are standardized using a StandardScaler (Pedregosa et al., 2011).

Dataset	train	validation	test
Electricity	70%	10%	20%
Weather	70%	10%	20%
ETTh1	60%	20%	20%
ETTM1	60%	20%	20%
AG News	90%	0%	10%

Table 4: Dataset splitting for the benchmark datasets.

**Hyperparameter Settings.** The key hyperparameter settings are listed in Table 5. The prediction lengths are all given {96, 192, 336, 720}.

Dataset	n-blocks	$\kappa$ -list	batch-size	input-len	epochs
Electricity	1	[3]	128	96	25
Weather	1	[3]	128	96	25
ETTh1	1	[3]	32	96	25
ETTM1	1	[3]	32	96	25
AG News	1	[3]	32	-	25

Table 5: Key hyperparameter settings on the benchmark datasets for Fighter.

We also adopt the ReduceLRonPlateau (Hinton et al., 2012) as the learning rate scheduler, and the early stopping for the ETTh1 and ETTm1 datasets.

### B.2 SENSITIVITY ANALYSIS

We present attention heatmaps for the first 64 positions at four training moments (25, 50, 75, and 100 rounds) and four kappa settings, *i.e.*,  $\kappa \in \{1, 2, 3, 4\}$ . As  $\kappa$  increases, a position is more likely to activate the most relevant position, effectively expanding the receptive field. By adjusting  $\kappa$  in time series and text classification tasks, Fighter gradually forms a collaborative modeling of global and local information, improving its representation capabilities.

The attention heatmaps in Table 8 provide a detailed visualization of Fighter and the Transformer model, illustrating the impact of higher-order dynamic attention mechanism and the connection to graph-based feature aggregation. This paper reinterprets the Transformer encoder through the lens of graph convolution, where the attention distribution matrix acts as a dynamic adjacency matrix. In this context, Fighter extends the Transformer by incorporating higher-order compositions of the attention distribution matrix, controlled by the parameter  $\kappa$ , which governs the power to which the attention distributions are raised. This flexible mechanism enables Fighter to perform multi-hop graph aggregation, capturing temporal dependencies across multiple scales. The heatmaps vividly demonstrate how  $\kappa$  transforms the Fighter’s attention patterns, evolving from localized dependencies to capturing broader, global dependencies.

It can be observed that the Transformer’s attention heatmap is constrained to local dependencies, as reflected in its strong diagonal focus. While this is effective for short-range temporal modeling, it limits the Transformer’s capacity to capture long-range dependencies and multi-scale relationships.

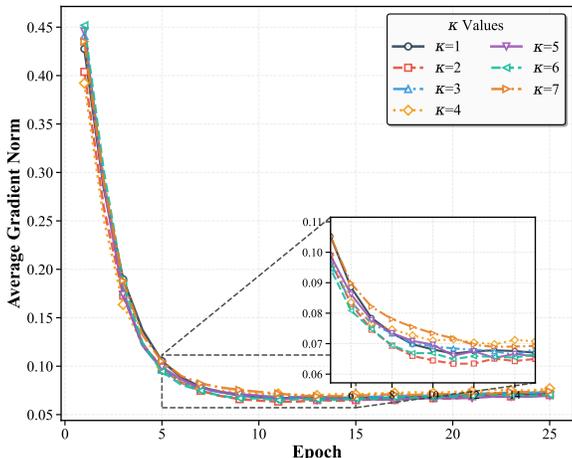
1188 Fighter, on the other hand, introduces higher-order attention distribution matrix to aggregate infor-  
 1189 mation across multiple hops in the time series. At  $\kappa = 2$ , Fighter’s attention heatmap resembles  
 1190 the Transformer’s, emphasizing local interactions. However, as  $\kappa$  increases, the Fighter’s attention  
 1191 becomes more distributed and structured, capturing long-range dependencies and revealing richer  
 1192 temporal patterns. This behavior aligns with the theoretical interpretation of Fighter as performing  
 1193 higher-order graph convolution, where increasing  $\kappa$  corresponds to aggregating information from  
 1194 wider neighborhoods in the graph.

1195 The incorporation of higher-order attention mechanisms in Fighter provides several key advantages.  
 1196 By exponentiating the attention distribution matrix, Fighter flexibly emphasizes stronger relationships  
 1197 and suppresses weaker ones, effectively controlling the breadth of the graph aggregation process.  
 1198 This mechanism not only improves the model’s ability to capture complex temporal dependencies but  
 1199 also provides an interpretable representation of these dependencies as graph edges. The parameter  
 1200  $\kappa$  acts as a tunable knob, balancing local and global attention patterns, thereby enabling Fighter to  
 1201 adapt to the specific demands of the task. This flexibility, combined with the removal of redundant  
 1202 linear projections, results in a streamlined architecture that is both computationally efficient and  
 1203 interpretable.

1204 Fighter bridges the gap between Transformers and GCNs, offering a unified framework for time  
 1205 series modeling. By incorporating higher-order graph aggregation through  $\kappa$ , Fighter achieves a  
 1206 multi-scale understanding of temporal dependencies, naturally expressed as graph relationships. The  
 1207 heatmaps corroborate Fighter’s ability to capture both local and global dependencies, providing  
 1208 a mechanistically interpretable view of its predictions. These insights, validated by competitive  
 1209 performance on standard forecasting benchmarks, highlight Fighter’s potential as a powerful and  
 1210 interpretable architecture for time series modeling and beyond.

1211 C ADDITIONAL EXPERIMENT

1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228



1229 Figure 6: The training average gradient norm curves across different  $\kappa$ .

1230  
1231  
1232 In Table 6, we record the GPU memory usage comparison for Fighter ( $\kappa = 3$ ) and the baselines on  
 1233 the Weather dataset when the prediction length is 96. Fighter ( $\kappa = 3$ ) implements the multi-hop  
 1234 graph aggregation, which enables long-range information to propagate naturally by higher-order  
 1235 compositions of the attention distribution matrix. As expected, it consumes relatively more GPU  
 1236 memory than the baselines but obtains superior performance. Furthermore, when  $\kappa = 1$ , Fighter  
 1237 exactly recovers a simplified version of the vanilla Transformer in which (i) the attention distribution  
 1238 matrix effectively collapses to the identity matrix, and (ii) the redundant second linear projection of  
 1239 the feed-forward network is removed (as justified by our GCN equivalence). Consequently, Fighter  
 1240 ( $\kappa = 1$ ) consumes noticeably less GPU memory than the vanilla Transformer. This confirms that  
 1241 the extra memory cost of Fighter is incurred only by the multi-hop aggregation branches ( $\kappa \geq 1$ )  
 and that, even at  $\kappa = 1$ , our streamlined design already brings a small but non-trivial memory  
 saving without any performance degradation. Notably, the overhead is linear in  $\kappa$ , as presented

in Table 7. Although Fighter( $\kappa = 3$ ) introduces additional computational branches, resulting in a moderate increase in GPU memory usage, the overall consumption remains lightweight and well within practical limits. This overhead is justified by the substantial performance gains Fighter delivers compared with conventional Transformer-based models.

During training, we also record the average gradient norm comparison of Fighter across different  $\kappa$  in Figure 6. As demonstrated, the average gradient norm of Fighter across varying  $\kappa$  exhibits a vivid trend toward optimization stability: it starts at a relatively high value in the early training epochs and gradually decreases to a small and stable range as training proceeds. Importantly, the introduction of  $\kappa$  does not disrupt the stability of the optimization dynamics. The gradient norms of Fighter( $\kappa = 3$ ) follows a similar decay pattern and maintain comparable fluctuation levels to those corresponding to other  $\kappa$  values. This indicates that the  $\kappa$  does not introduce gradient explosion or vanishing phenomena, meanwhile, does not lead to abnormal oscillations or instability during training.

Table 6: GPU memory usage for Fighter ( $\kappa = 3$ ) and different baselines on the Weather dataset (the prediction length is 96).

Model	Autoformer	Informer	Reformer	Transformer	Fighter( $\kappa = 3$ )
<b>GPU Mem (MB)</b>	108.0	120.2	105.6	109.1	157.3

Table 7: GPU memory usage for Fighter ( $\kappa \in [1, 7]$ ) on the Weather dataset (the prediction length is 96).

Model	$\kappa = 1$	$\kappa = 2$	$\kappa = 3$	$\kappa = 4$	$\kappa = 5$	$\kappa = 6$	$\kappa = 7$
<b>GPU Mem (MB)</b>	94.5	125.4	157.3	188.8	233.2	278.7	337.1

Table 8: Visualization of attention heatmaps for the first 64 keys, comparing **Fighter** and Transformer. Rows indicate the models; columns represent training epochs (25, 50, 75, 100).

