

MSAGAT-Net: Multi-Scale Adaptive Graph Attention Network for Spatiotemporal Forecasting

Architecture Diagrams

April 20, 2025

1 Network Architecture Overview

2 Key Architectural Components

2.1 Depthwise Separable Convolution for Temporal Feature Extraction

2.2 Efficient Adaptive Graph Attention Module

2.3 Dilated Multi-Scale Temporal Module

2.4 Progressive Prediction Module

3 Tensor Flow Throughout Network

4 Mathematical Formulation

The MSAGAT-Net model can be described mathematically as follows:

$$\mathbf{X} \in \mathbb{R}^{B \times T \times N} \quad (\text{Input time series}) \quad (1)$$

$$\mathbf{A} \in \mathbb{R}^{N \times N} \quad (\text{Adjacency matrix}) \quad (2)$$

$$\mathbf{F} = \text{TemporalFeatureExtraction}(\mathbf{X}) \quad (\text{Temporal features}) \quad (3)$$

$$\mathbf{H} = \text{FeatureProcessing}(\mathbf{F}) \quad (\text{Processed features}) \quad (4)$$

$$\mathbf{H}' = \text{EfficientGraphAttention}(\mathbf{H}, \mathbf{A}) \quad (\text{Spatial features}) \quad (5)$$

$$\mathbf{H}'' = \text{MultiScaleTemporal}(\mathbf{H}') \quad (\text{Multi-scale features}) \quad (6)$$

$$\hat{\mathbf{Y}} = \text{ProgressivePrediction}(\mathbf{H}'', \mathbf{X}_T) \quad (\text{Predictions}) \quad (7)$$

$$\mathcal{L} = \text{MSE}(\mathbf{Y}, \hat{\mathbf{Y}}) + \lambda \|\boldsymbol{\alpha}\|_1 \quad (\text{Loss function}) \quad (8)$$

where \mathbf{X}_T represents the last observed values and $\lambda \|\boldsymbol{\alpha}\|_1$ is the attention regularization term that encourages sparse attention patterns for better interpretability.

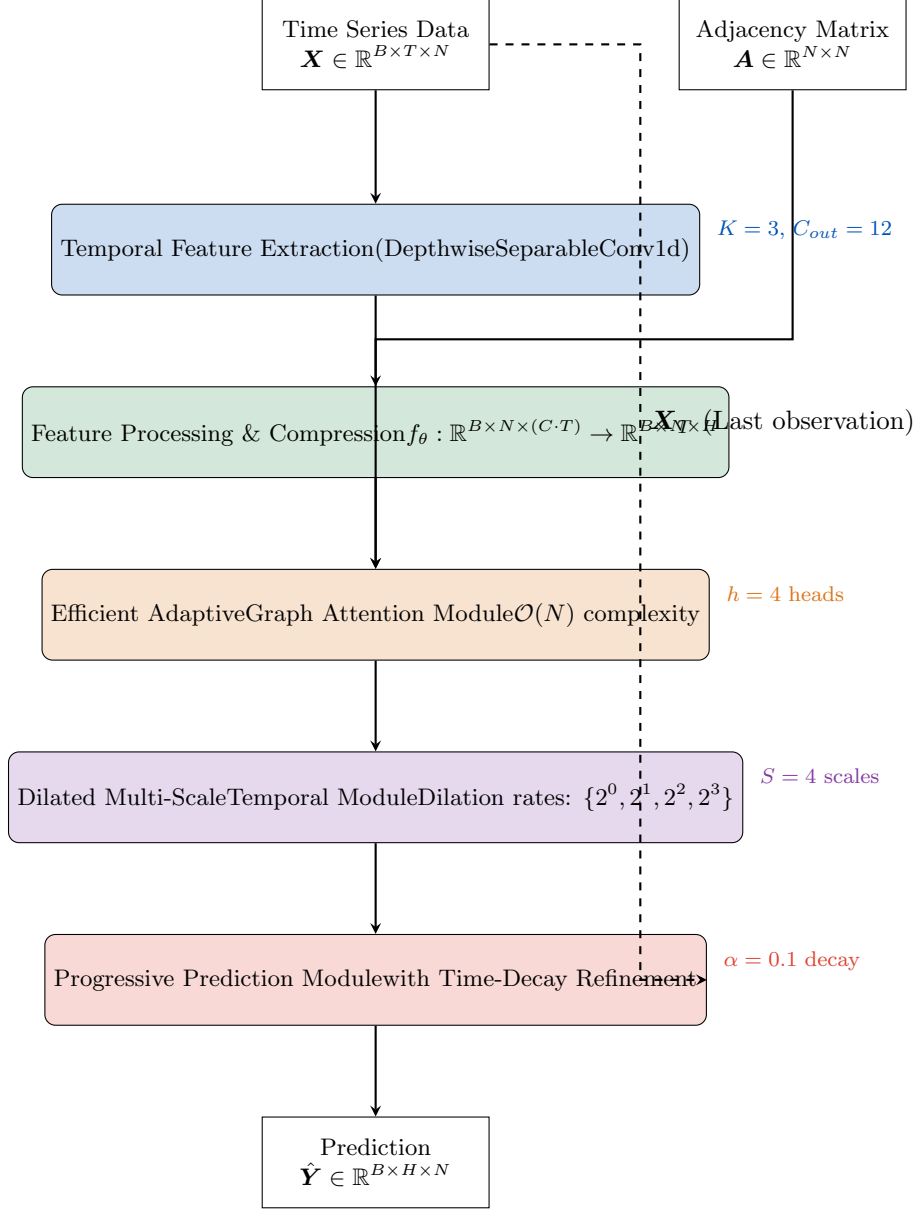
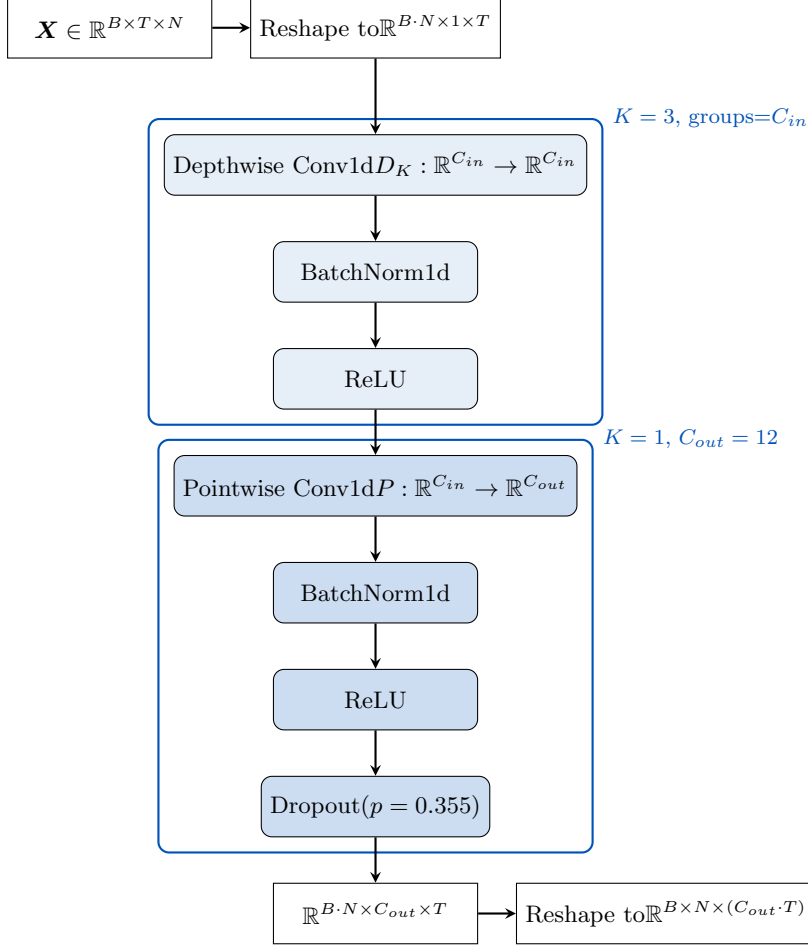


Figure 1: Overall architecture of MSAGAT-Net for spatiotemporal forecasting. The model processes time series data $\mathbf{X} \in \mathbb{R}^{B \times T \times N}$ and an adjacency matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ through specialized modules for temporal and spatial feature extraction. B represents batch size, T is the input window size, N is the number of nodes, H is the hidden dimension, and the output $\hat{\mathbf{Y}}$ provides predictions for the next H time steps.



$$D_K(X)_{i,j,t} = \sum_{k=1}^K W_{i,k}^d \cdot X_{i,j,t+k-\frac{K+1}{2}}$$

$$P(X)_{i,j,t} = \sum_{c=1}^{C_{in}} W_{i,j,c}^p \cdot X_{c,j,t}$$

Parameters: $\mathcal{O}(K \cdot C_{in} + C_{in} \cdot C_{out})$

Figure 2: Depthwise Separable Convolution for temporal feature extraction. This factorized convolution operation first applies a depthwise convolution D_K with kernel size $K = 3$ that processes each channel independently, followed by a pointwise convolution P that mixes the channels. This factorization reduces the parameter count from $\mathcal{O}(K \cdot C_{in} \cdot C_{out})$ to $\mathcal{O}(K \cdot C_{in} + C_{in} \cdot C_{out})$ while maintaining strong representational capacity.

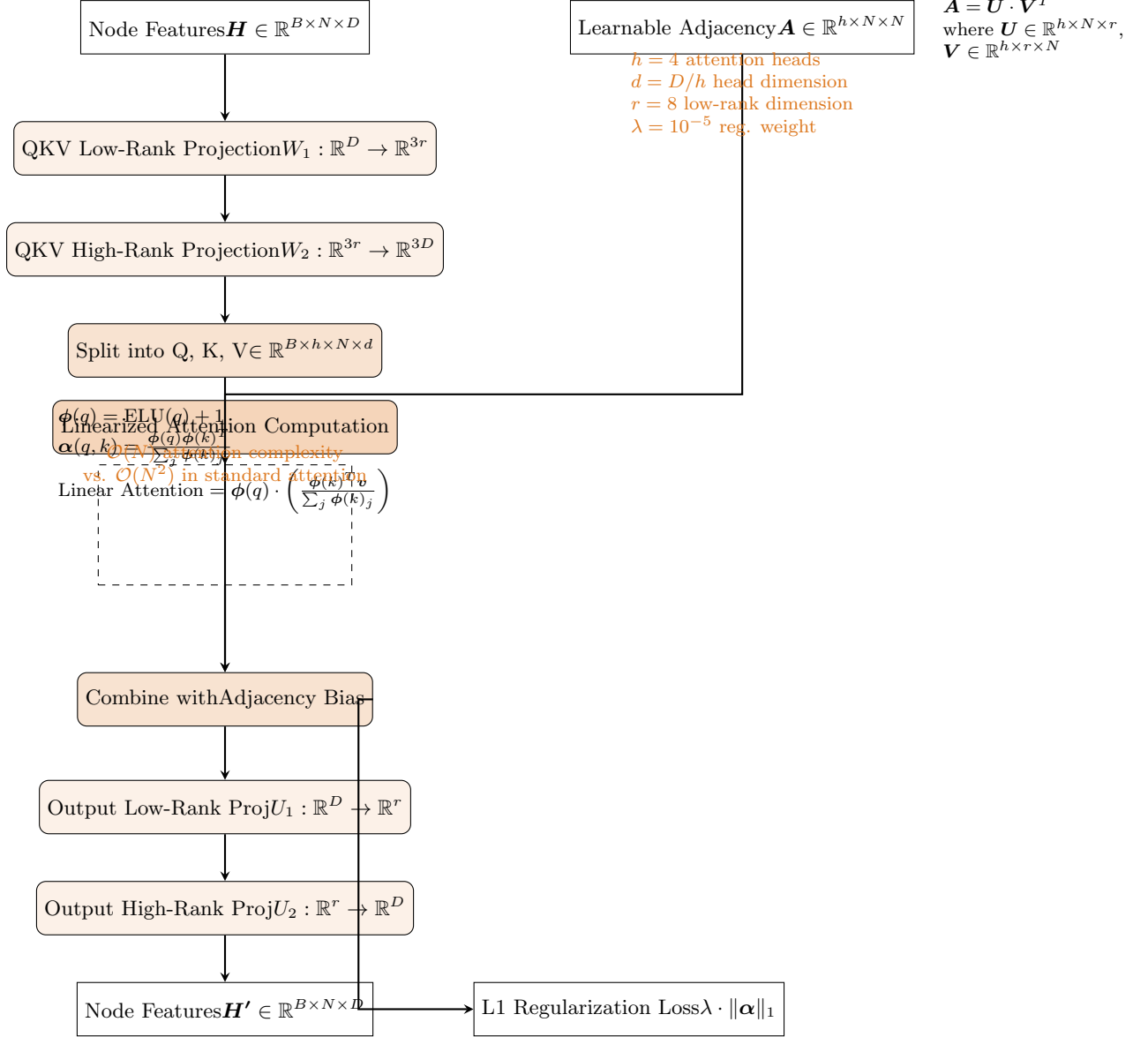


Figure 3: Efficient Adaptive Graph Attention Module with $\mathcal{O}(N)$ complexity and low-rank factorization. This module implements linear attention using the ELU+1 kernel trick, which transforms the quadratic complexity of standard attention to linear. The learnable adjacency matrix is factorized as $\mathbf{A} = \mathbf{U} \cdot \mathbf{V}^T$ to reduce memory requirements from $\mathcal{O}(N^2)$ to $\mathcal{O}(Nr)$. All projection matrices are decomposed using low-rank factorization for parameter efficiency. The L1 regularization encourages sparse attention patterns for better interpretability.

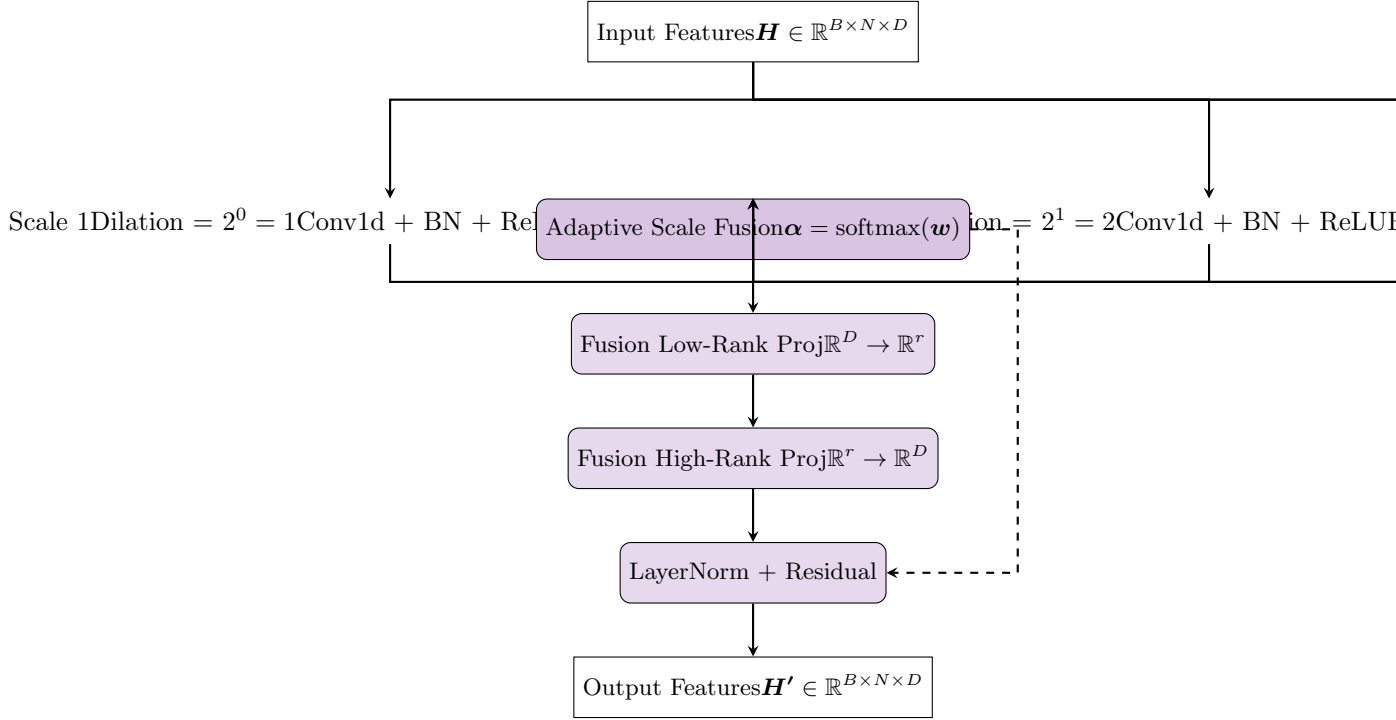
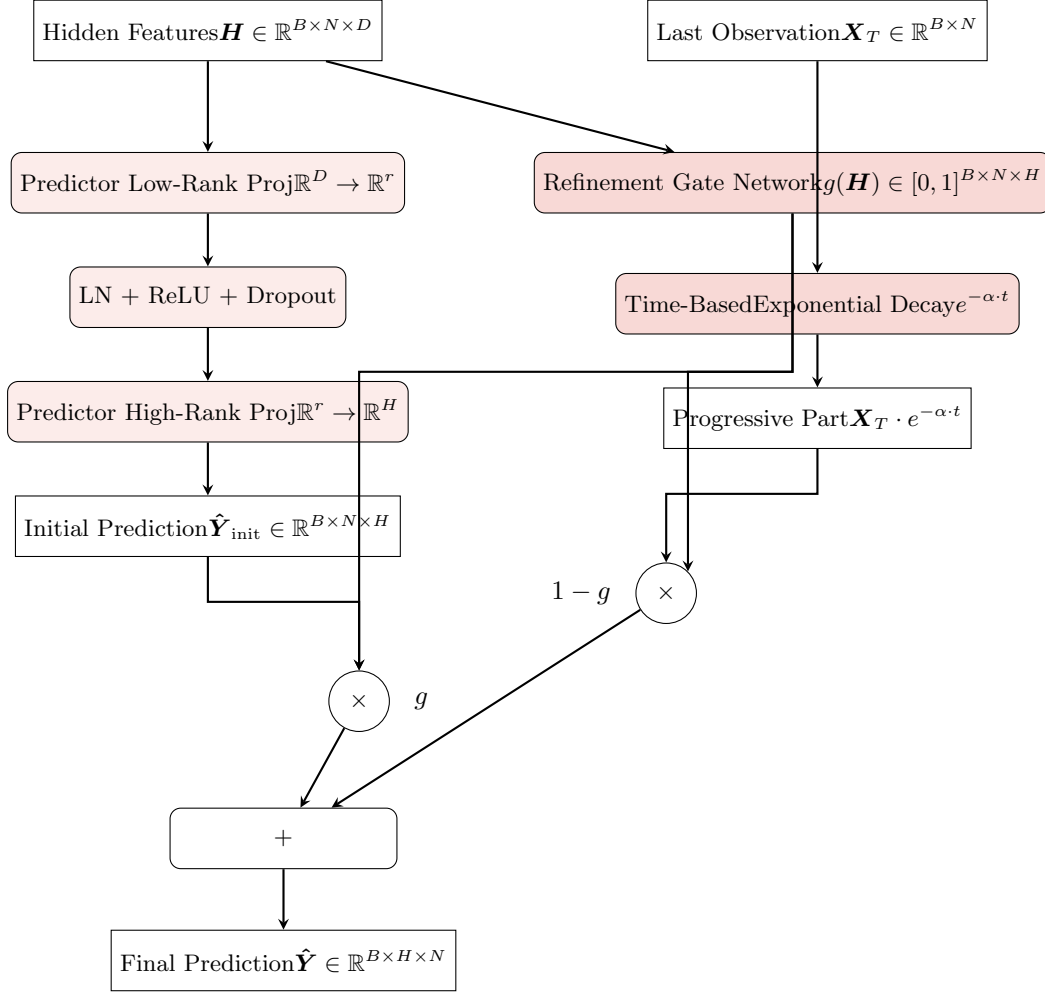


Figure 4: Dilated Multi-Scale Temporal Module with adaptive fusion of features from different time scales. This module captures temporal patterns at multiple resolutions by using parallel dilated convolutions with exponentially increasing dilation rates (1, 2, 4, 8). This creates an effective receptive field ranging from 3 to 17 time steps without requiring deep stacking of layers. The features from different scales are combined using learned weights that are dynamically adapted during training. A residual connection helps maintain gradient flow.



$$g(\mathbf{H}) = \sigma(W_2 \cdot \text{ReLU}(W_1 \mathbf{H}))$$

$$\mathbf{p}_t = \mathbf{X}_T \cdot e^{-\alpha t}$$

$$\hat{\mathbf{Y}} = g \cdot \hat{\mathbf{Y}}_{\text{init}} + (1 - g) \cdot \mathbf{p}$$

where $\alpha = 0.1$ is decay rate,
 $t \in \{1, 2, \dots, H\}$ is horizon step

Figure 5: Progressive Prediction Module with adaptive gating and time-decay refinement. This module generates initial multi-step predictions that are then adaptively combined with a persistence forecast derived from the last observation \mathbf{X}_T . The gating network dynamically determines how much to rely on the model’s prediction versus the persistence forecast for each node and time step. The persistence forecast is adjusted with an exponential decay factor $e^{-\alpha t}$ to reflect decreasing confidence in persistent patterns over longer horizons. This approach helps stabilize predictions, especially for shorter horizons.

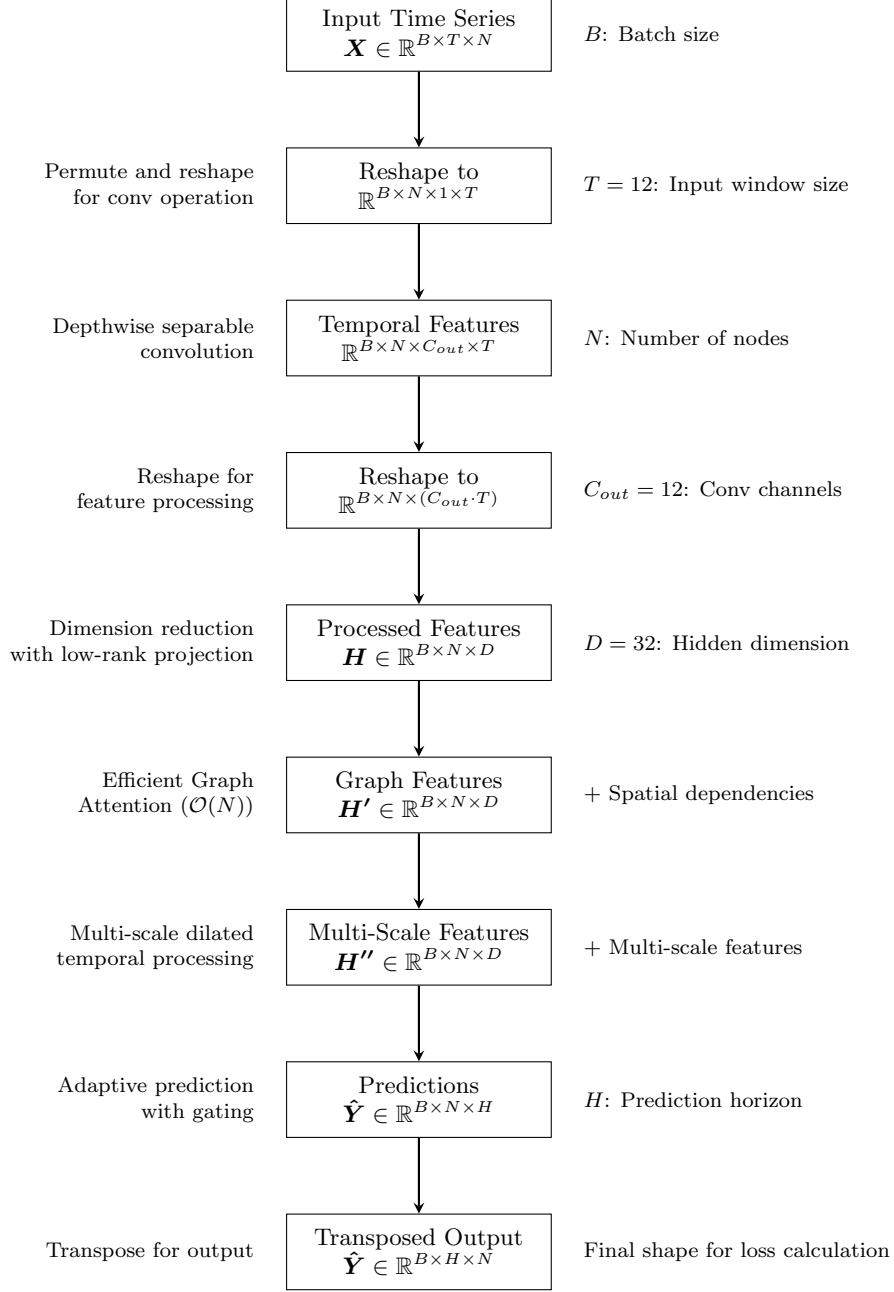


Figure 6: Detailed tensor flow through the MSAGAT-Net architecture showing dimensions and transformations at each stage. The network processes input time series data with shape $[B, T, N]$ through a sequence of specialized modules that extract and combine temporal and spatial features. Low-rank projections are used throughout the network to reduce parameters while maintaining model capacity. The final output tensor $[B, H, N]$ provides predictions for each node across the forecast horizon.