



Backpropagation Computation for Training Graph Attention Networks

Joe Gould¹ · Keshab K. Parhi¹

Received: 3 June 2023 / Revised: 27 September 2023 / Accepted: 28 September 2023 / Published online: 16 October 2023
© The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2023

Abstract

Graph Neural Networks (GNNs) are a form of deep learning that have found use for a variety of problems, including the modeling of drug interactions, time-series analysis, and traffic prediction. They represent the problem using non-Euclidian graphs, allowing for a high degree of versatility, and are able to learn complex relationships by iteratively aggregating more contextual information from neighbors that are farther away. Inspired by its power in transformers, Graph Attention Networks (GATs) incorporate an attention mechanism on top of graph aggregation. GATs are considered the state of the art due to their superior performance. To learn the best parameters for a given graph problem, GATs use traditional backpropagation to compute weight updates. To the best of our knowledge, these updates are calculated in software, and closed-form equations describing their calculation for GATs aren't well known. This paper derives closed-form equations for backpropagation in GATs using matrix notation. These equations can form the basis for design of hardware accelerators for training GATs.

Keywords Neural network training · Backpropagation · Gradient computation · Graph attention networks

1 Introduction

Many complex problems can be represented by graphs, and graph neural networks (GNNs) have been used to address these issues. GNNs have been successfully applied to applications such as drug interactions [1, 2], time-series analysis [3], traffic prediction [4], neurological damage detection [5], and others. One of the best performing GNN architectures are Graph Attention Networks (GATs) [6]. These networks build upon previous models by combining the statically weighted graph convolution with an input-dependent attention mechanism to achieve higher expressive power.

Although GAT networks are powerful, they are computationally expensive, requiring processing for an attention coefficient for every edge in the graph. While the forward pass, or inference, of GAT network processing is well understood, to the best of our knowledge, closed-form expressions for backpropagation for training GAT networks have not been presented before. The main contribution of this paper

is the derivation of these equations using matrix notation for both the original GAT and a similar derivative work, GATv2 [7]. These equations can enable design of accelerators for training GATs by exploiting techniques such as gradient interleaving [8].

This paper is organized as follows. Section 2 presents a brief review of GNNs and GATs. Section 3 presents the forward pass equations associated with the inference from GATs in matrix form. In Section 4, we derive the backpropagation equations for training the GATs in matrix form. Section 5 concludes the paper. Large figures that may be helpful alongside the text are found in Section of [Appendix](#).

2 Background

2.1 Graph Neural Networks

GNNs are a broad class of architectures that developed from early work with Recurrent Neural Networks [9, 10] that use deep neural networks to perform different graph-based tasks [11, 12]. The most relevant type of GNN for understanding GATs is Graph Convolutional Network (GCN) [13], which many GNN derivatives are based on. GCN and GAT operate by performing spatial graph convolution and traditional machine learning transformations on a graph to extract meaningful information over a wider

✉ Joe Gould
gould146@umn.edu

✉ Keshab K. Parhi
parhi@umn.edu

¹ Department of Electrical and Computer Engineering,
University of Minnesota, Minneapolis 55455, MN, USA

region. GCNs and GAT networks consist of layers, with each performing the convolution steps and applying a non-linear activation before outputting into the next layer.

Spatial graph convolution can be described as a form of message passing [14] on graphs. Nodes have features which are used to produce a message, the messages are transformed by and passed along edges, and the messages are gathered at edge destinations and are used to update the destination node states. This description is general, and GCN simplifies computation by limiting the convolution operation. In this modified convolution operation, edges have scalar weights which are multiplied with source node features. These weighted features are summed at the destination node, and the resultant sum is linearly transformed with learnable parameters. This convolution can be written using vector notation,

$$\mathbf{z}_i = \sum_{j \in \mathcal{N}(i)} \mathbf{A}_{ij} \cdot \mathbf{h}_j \times \mathbf{W} \tag{1}$$

where $\mathcal{N}(i)$ is the neighborhood of node i , \mathbf{A}_{ij} is the scalar edge weight from node j to node i , \mathbf{W} is the learnable weight matrix, and \mathbf{h}_i and \mathbf{z}_i are row vectors and the pre- and post-convolution features of node i . A similar operation is used in GAT.

As can be seen from the description above, the convolution is able to propagate information across the structure of a graph. This has the effect of allowing subsequent node representations to incorporate a wider neighborhood into their extracted features, and is the main contributor to the expressive power of GCNs and GATs. It is analogous to the widening receptive field of pixels in traditional Convolutional Neural Networks [11]. This is illustrated in Fig. 1 for a 2-layer network on a graph of 8 nodes.

A number of hardware accelerators exist for GNN processing. The forward pass is well understood, and

architectures which can support GAT inference exist [15–17]. However, architectures which support training are usually geared towards training GCN [18–20] or accelerating software solutions [21, 22] for training GAT [23]. While these methods work, it is possible a more powerful accelerator could be created by specifically targeting the computations for GAT training operations.

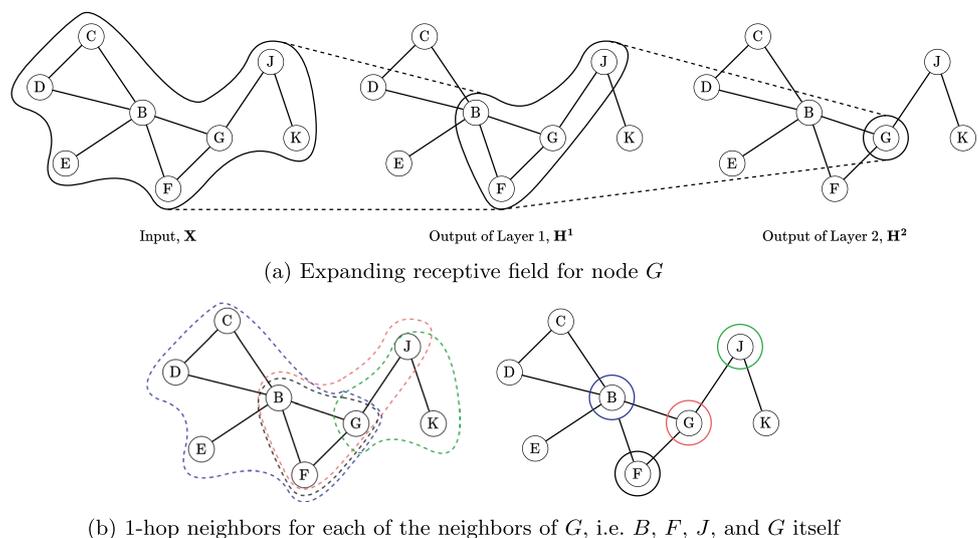
2.2 Neural Network Attention

Attention is a concept in machine learning that has been extensively used in applications such as machine translation [24, 25] and machine vision [26], and is the primary mechanism that transformers use for their generalized effectiveness [27]. It is based on the intuition that prioritizing the important features of an input will allow a network to better extract information from it. A network layer typically implements attention by performing a nonlinear transformation on the inputs with some learned attention parameters, and then using a softmax operation for each of the outputs. This has the effect of allowing the network to learn input features that are indicative of importance and then raising their corresponding weights to focus on them.

3 GAT Forward Pass

GAT networks operate similar to GCNs, but implement an attention mechanism to adjust edge coefficients based on node features. Like other neural network types, GATs use stacked layers to extract information from inputs and generate useful output representations. Each layer has its own set of learnable parameters, which are updated during training. Taking from existing literature on attention [27], GAT splits processing steps within each layer across multiple attention

Figure 1 Example GNN with 2 layers on a graph of 8 nodes.



heads, which have their own attention weights during the aggregation step. Baseline GAT consists of the following processing within a layer, for each attention head:

1. Node features are linearly transformed by learnable weights to form combined features.
2. Each node’s combined features are transformed into source and destination coefficients.
3. For every edge, a scalar value is computed using its source and destination coefficients. The scalar values of the edges sharing a common destination node are input to a softmax function to generate the attention coefficients of the edges.
4. GCN-like aggregation is performed using the original combined features and attention coefficients. A nonlinear activation is applied, and the result is concatenated with the results from the other attention heads.

There are currently two forms of GATs in use, the original GAT and GATv2. Hereafter, the original GAT will be referred to as GATv1, and GAT will be used to refer to the general architecture of these networks. GATv2 was proposed because of issues in GATv1’s attention mechanism, which made it unable to effectively differentiate between certain types of inputs. However, because GATv1 has been in use for longer than GATv2, we present derivations for both versions. We begin by constructing equations for GAT inference in matrix form. We first define the relevant parameters in Table 1, and then follow each step of layer processing as described above.

We use the superscripts l to mean pertaining to layer l , and k, l to mean pertaining to attention head k of layer l . Non-scalar variables are denoted in bold. Vectors are assumed to be columns, unless they are taken from rows of matrices, and use lower case. When taking sub-elements of a multidimensional variable, we use subscripts to denote the row or column being taken from, and keep the same bolding and case as the variable the sub-elements are taken from. We consider graphs with directed edges, and do not require the adjacency matrix to be symmetric.

3.1 Linear Transformation of Node Features

The layer begins by linearly transforming the output of the previous layer to a new feature, not necessarily of the same size, for the current layer. GATv1 uses a single learnable weight matrix common to every node, while GATv2 has two matrices, with the result of the combination used differently depending on if a node is a source or destination of an edge.

For GATv1, this step can be written as

$$\mathbf{X}^{k,l} = \mathbf{H}^{l-1} \times \mathbf{W}^{k,l}. \tag{2}$$

GATv2 has two combined matrices, which are

$$\begin{aligned} \mathbf{X}_{src}^{k,l} &= \mathbf{H}^{l-1} \times \mathbf{W}_{src}^{k,l}, \\ \mathbf{X}_{dst}^{k,l} &= \mathbf{H}^{l-1} \times \mathbf{W}_{dst}^{k,l}. \end{aligned} \tag{3}$$

3.2 Edge Coefficient Computation

Using the combined features, GAT generates a set of source and destination coefficients that are used for computing attention coefficients. In GATv1, these coefficients are scalars obtained from performing a dot product between the learnable attention weights and the combined feature vector. In GATv2, the coefficients are vectors, and are just the combined features generated in step 1. GATv2 introduces the attention weights in the next step.

For GATv1, the scalar coefficients for node i are

$$\begin{aligned} c_{i,V1,src}^{k,l} &= \mathbf{a}_{src}^{k,lT} \times \mathbf{X}_i^{k,l} = \mathbf{a}_{src}^{k,lT} \times \mathbf{W}_{src}^{k,lT} \times \mathbf{h}_i^{k,l}, \\ c_{i,V1,dst}^{k,l} &= \mathbf{a}_{dst}^{k,lT} \times \mathbf{X}_i^{k,l} = \mathbf{a}_{dst}^{k,lT} \times \mathbf{W}_{dst}^{k,lT} \times \mathbf{h}_i^{k,l}, \\ \mathbf{c}_{V1,src}^{k,l} &= \left\| \left\| c_{i,V1,src}^{k,l} \right\|_{i \in \mathcal{N}} \right\|, \\ \mathbf{c}_{V1,dst}^{k,l} &= \left\| \left\| c_{i,V1,dst}^{k,l} \right\|_{i \in \mathcal{N}} \right\|. \end{aligned} \tag{4}$$

The vector coefficients for GATv2 are

$$\begin{aligned} \mathbf{c}_{i,V2,src}^{k,l} &= \mathbf{X}_{i,src}^{k,l} = \mathbf{W}_{src}^{k,lT} \times \mathbf{h}_i^{k,l}, \\ \mathbf{c}_{i,V2,dst}^{k,l} &= \mathbf{X}_{i,dst}^{k,l} = \mathbf{W}_{dst}^{k,lT} \times \mathbf{h}_i^{k,l}, \\ \mathbf{C}_{V2,src}^{k,l} &= \left\| \left\| \mathbf{c}_{i,V2,src}^{k,l} \right\|_{i \in \mathcal{N}} \right\| = \mathbf{X}_{src}^{k,lT}, \\ \mathbf{C}_{V2,dst}^{k,l} &= \left\| \left\| \mathbf{c}_{i,V2,dst}^{k,l} \right\|_{i \in \mathcal{N}} \right\| = \mathbf{X}_{dst}^{k,lT}. \end{aligned} \tag{5}$$

3.3 Attention Coefficient Computation

Next, every edge in the graph has a scalar value computed. This is done by first adding the source and destination coefficients from the source and destination nodes of each edge, respectively. After this sum is computed, it is put through a nonlinear leaky ReLU (LReLU) operation. In GATv1, the result of this is the pre-softmax normalization attention coefficient, while in GATv2 it is a vector that is used in an inner product with the learnable attention coefficients. The result of the GATv2 inner product is its pre-softmax normalization attention coefficient. For GATv1, this can be written as

$$e_{ij}^{k,l} = \begin{cases} \sigma_L \left(c_{i,V1,dst}^{k,l} + c_{j,V1,src}^{k,l} \right) & \mathbf{A}_{ij} \neq 0, \\ -\infty & \mathbf{A}_{ij} = 0 \end{cases}, \tag{6}$$

Table 1 Notations for GAT Inference.

Notation		Description
N		Number of nodes in the graph
L		Layers within network
l		Index of a layer within network, $\in [1, L]$
K^l		Number of attention heads in layer l
k		Index of an attention head within layer, $\in [1, K^l]$
d^l		Output feature dimension for an attention head in layer l
k^0		Defined as 1
d^0		Defined as the feature dimension of input data
$\sigma_L, \sigma_S, \sigma_E$		LReLU, softmax, and ELU nonlinear functions, respectively
$\parallel \dots$		Concatenation of elements horizontally
\oplus		Defined in Eqs. (6) and (7), shown visually in Fig. 4 subfigure c
Notation	Size	Description
\mathbf{A}	$\mathbb{R}^{N \times N}$	Adjacency matrix of graph. The i th row corresponds to the incoming edges to node i , and the j th column corresponds to the outgoing from node j . All entries are either 0 or 1.
\mathbf{h}_i^l	$\mathbb{R}^{(K^l \cdot d^l)}$	Output features of node i from layer l .
\mathbf{H}^l	$\mathbb{R}^{N \times (K^l \cdot d^l)}$	Output features from layer l . The i th row corresponds to the feature vector of node i .
$\mathbf{W}^{k,l}, \mathbf{W}_{src}^{k,l}, \mathbf{W}_{dst}^{k,l}$	$\mathbb{R}^{(K^{l-1} \cdot d^{l-1}) \times d^l}$	Learnable weight matrices for attention head k in layer l . Separate among source/destination in GATv2, common in GATv1.
$\mathbf{X}^{k,l}, \mathbf{X}_{src}^{k,l}, \mathbf{X}_{dst}^{k,l}$	$\mathbb{R}^{N \times d^l}$	Combined feature matrices for attention head k in layer l , from their respective weight matrices.
$\mathbf{a}^{k,l}, \mathbf{a}_{src}^{k,l}, \mathbf{a}_{dst}^{k,l}$	\mathbb{R}^{d^l}	Learnable source/destination attention weights for attention head k in layer l . Separate among source/destination in GATv1, common in GATv2.
$c_{i,V1,src}^{k,l}, c_{i,V1,dst}^{k,l}$	\mathbb{R}	GATv1: Source/destination scalar coefficients for node i in attention head k of layer l .
$\mathbf{c}_{V1,src}^{k,l}, \mathbf{c}_{V1,dst}^{k,l}$	\mathbb{R}^N	and collected into a vector.
$c_{i,V2,src}^{k,l}, c_{i,V2,dst}^{k,l}$	\mathbb{R}^{d^l}	GATv2: Source/destination vector coefficients for node i in attention head k of layer l .
$\mathbf{C}_{V2,src}^{k,l}, \mathbf{C}_{V2,dst}^{k,l}$	$\mathbb{R}^{N \times d^l}$	and collected into a matrix.
$e_{ij}^{k,l}, \alpha_{ij}^{k,l}$	\mathbb{R}	Pre- and post-softmax normalization attention coefficients for the edge outgoing from node j to node i .
$\mathbf{L}^{k,l}, \mathbf{S}^{k,l}$	$\mathbb{R}^{N \times N}$	Pre- and post-softmax attention matrices. Entry (i, j) is the corresponding weight from node j to node i .
$\mathbf{Z}^{k,l}$	$\mathbb{R}^{N \times (K^l \cdot d^l)}$	Pre-activation output features for layer l .

and in GATv2 as

$$e_{ij}^{k,l} = \begin{cases} \mathbf{a}^{k,l} \times \sigma_L \left(\mathbf{c}_{i,V2,dst}^{k,l} + \mathbf{c}_{j,V2,src}^{k,l} \right) & \mathbf{A}_{ij} \neq 0 \\ -\infty & \mathbf{A}_{ij} = 0 \end{cases}. \quad (7)$$

We denote the addition operation with \oplus . In the case of GATv1, $\mathbf{c}_{V1,src}^{k,l}$ and $\mathbf{c}_{V1,dst}^{k,l}$ are vectors, and the operation generates a matrix by performing an outer sum. In GATv2, $\mathbf{C}_{V2,dst}^{k,l}$ and $\mathbf{C}_{V2,src}^{k,l}$ are matrices, and the operation generates a 3-dimensional tensor by performing the outer sum on the columns of the matrices. Using this notation, we can write the resulting attention matrices as

$$\mathbf{L}^{k,l} = \sigma_L \left(\mathbf{c}_{V1,dst}^{k,l} \oplus \mathbf{c}_{V1,src}^{k,l} \right) \odot \mathbf{A} \quad (8)$$

for GATv1, and for GATv2 as

$$\mathbf{L}^{k,l} = \left(\mathbf{a}^{k,l} \times \sigma_L \left(\mathbf{C}_{V2,dst}^{k,l} \oplus \mathbf{C}_{V2,src}^{k,l} \right) \right) \odot \mathbf{A}. \quad (9)$$

Finally, in both GATv1 and GATv2, these coefficients are put through a softmax, grouped by common destinations, to normalize them to the same dynamic range. These matrices have an identical nonzero pattern to the adjacency matrix \mathbf{A} . This results in the \mathbf{L} and \mathbf{S} matrices that have their (i, j) th entries be the pre- and post-softmax attention coefficient for the edge outgoing from node j to node i as

$$\mathbf{S}_i^{k,l} = \sigma_S \left(\mathbf{L}_i^{k,l} \right). \quad (10)$$

3.4 GCN-like Aggregation and Activation

The final step in the forward pass is to propagate the combined node features using the final attention coefficient on each edge. This propagation is identical to GCN, with the adjacency matrix replaced with the post-softmax attention matrix. By performing matrix multiplication in this way, the

combined features of all of the in-neighbors of each node are summed together, weighed by their attention. In GATv1, this is given by

$$\mathbf{Z}^{k,l} = \mathbf{S}^{k,l} \times \mathbf{X}^{k,l}. \tag{11}$$

This is identical to GATv2, with the features combined with the source weights matrix in the aggregation, written as

$$\mathbf{Z}^{k,l} = \mathbf{S}^{k,l} \times \mathbf{X}_{\text{src}}^{k,l}. \tag{12}$$

Finally, the combined and aggregated features are put through a nonlinear activation. We use σ_E to represent the exponential linear unit (ELU) activation function. The resulting features for each attention head are concatenated together to form the output of the layer:

$$\mathbf{H}^{l+1} = \parallel_{k=1}^{K^l} \sigma_E(\mathbf{Z}^{k,l}). \tag{13}$$

In the case of the output layer, multiple attention heads have their results averaged instead of being concatenated, and a nonlinear operation like softmax is applied on the resulting node features. This can be written as

$$\mathbf{H}_i^L = \sigma_S \left(\frac{1}{K^L} \cdot \sum_{k=1}^{K^L} \mathbf{Z}_i^{k,L} \right). \tag{14}$$

We summarize the forward pass in a data flow graph (DFG) in Fig. 2. An illustrative example of each step for GATv1 and GATv2 in the forward pass is shown for a single layer’s attention head in Fig. 4. The example assumes an input feature size of 5 and an output feature size of the layer’s attention heads of 4.

4 GAT Backward Pass

We now derive the equations for the backpropagation for GAT’s learnable parameters, i.e., equations for gradients of the loss with respect to weight matrices $(\mathbf{W}^{k,l}, \mathbf{W}_{\text{src}}^{k,l}, \mathbf{W}_{\text{dst}}^{k,l})$, attention weights $(\mathbf{a}_{\text{src}}^{k,l}, \mathbf{a}_{\text{dst}}^{k,l}, \mathbf{a}^{k,l})$, and feature inputs (\mathbf{H}^l) for each layer. This will follow the same structure as the forward pass, starting with the input to the layers during backpropagation and computing the gradients associated with that layer. A summary of the notations for different variables relevant to the backward pass is provided in Table 2.

4.1 Gradients with Respect to Activation

We first compute the gradients of the loss with respect to the pre-activation outputs using the gradients input to the layer. The gradient can be computed as

$$\delta'^{k,l} = \sigma'_E(\mathbf{Z}^{k,l}) \odot \delta^{k,l}. \tag{15}$$

For the last layer in the network, a softmax activation function is used. Its derivative is well known, and the derivative of any vector through a softmax can be described by a symmetric gradient matrix,

$$\begin{aligned} \mathbf{u} &= \sigma_S(\mathbf{v}), \\ \frac{\partial u_i}{\partial v_j} &= \begin{cases} u_i - u_i^2 & i = j \\ -(u_i \cdot u_j) & i \neq j \end{cases}, \\ \frac{\partial \mathbf{u}}{\partial \mathbf{v}} &= \text{diag}(\mathbf{u}) - \mathbf{u} \times \mathbf{u}^T, \end{aligned}$$

where the (i, j) th entry in the $\frac{\partial \mathbf{u}}{\partial \mathbf{v}}$ matrix is the gradient contributed by the j th element of the input to the i th element of the output. In this case, \mathbf{v} and \mathbf{u} are the rows of the \mathbf{Z}^L and \mathbf{H}^L matrices. The gradient matrix described above can be computed for each node, and the gradient of the loss with respect to the pre-softmax features for that node can be computed by multiplying this matrix with the gradient of the loss with respect to the post-softmax features. This is shown as

$$\frac{\partial \mathbf{H}_i^L}{\partial \mathbf{Z}_i^L} = \text{diag}(\mathbf{H}_i^L) - \mathbf{H}_i^{L,T} \times \mathbf{H}_i^L, \tag{16}$$

$$\delta_i'^L = \frac{\partial \mathbf{H}_i^L}{\partial \mathbf{Z}_i^L} \times \delta_i^L. \tag{17}$$

For the rest of the backward pass, there is no difference between layers if we scale this pre-softmax gradient by K^L to account for the averaging step, written as

$$\delta'^{k,L} = \frac{1}{K^L} \cdot \delta'^L. \tag{18}$$

It is important to note that the matrices involved in this step are dense, unlike the matrix of attention coefficients. Because of the complexity of the softmax derivative and because it only needs to be computed for the output layer, δ'^L should be computed outside of a dedicated accelerator.

4.2 Gradients with Respect to Attention

We now look at the attention mechanism and how the gradient of the loss with respect to its inputs can be computed. The first step is to compute the gradient of the loss with respect post-softmax attention coefficients, through the aggregation, and then the pre-softmax edge coefficients, which has a form identical to the softmax derivative described before. This can be written as

$$\frac{\partial \mathcal{L}}{\partial \mathbf{S}^{k,l}} = \delta'^{k,l} \times \mathbf{X}^{k,l,T}, \tag{19}$$

$$\frac{\partial \mathbf{S}_i^{k,l}}{\partial \mathbf{L}_i^{k,l}} = \text{diag} \left(\mathbf{S}_i^{k,l} \right) - \mathbf{S}_i^{k,lT} \times \mathbf{S}_i^{k,l}, \quad (20)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{L}_i^{k,l}} = \frac{\partial \mathbf{S}_i^{k,l}}{\partial \mathbf{L}_i^{k,l}} \times \left(\frac{\partial \mathcal{L}}{\partial \mathbf{S}_i^{k,l}} \right)^T = \sigma'_s \left(\frac{\partial \mathcal{L}}{\partial \mathbf{S}_i^{k,l}} \right). \quad (21)$$

Because of the sparsity in $\mathbf{S}^{k,l}$, $\frac{\partial \mathbf{S}_i^{k,l}}{\partial \mathbf{L}_i^{k,l}}$ will only have nonzeros in rows/columns which correspond to the edges for node i . The complete loss with respect to the pre-softmax attentions, $\frac{\partial \mathcal{L}}{\partial \mathbf{L}^{k,l}}$, also has the same sparsity pattern as the adjacency matrix.

Now, because of the differences in GATv1 and GATv2, we separate the derivation for the rest of the attention backward pass for each of the two versions.

4.2.1 GATv1 Attention Gradient

The rest of the backward pass depends on the result of the LReLU that was used to compute the pre-softmax weight

matrix. Since the input to the LReLU was a scalar, we define a $\mathbf{C}^{k,l}$ matrix to hold values from the derivative of the LReLU, and then define the gradient with respect to its input as $\Delta^{k,l}$, which can be written as

$$\mathbf{C}'_{V1,ij}{}^{k,l} = \sigma'_L \left(c_{i,V1,dst}^{k,l} + c_{j,V1,src}^{k,l} \right), \quad (22)$$

$$\mathbf{C}'_{V1}{}^{k,l} = \sigma'_L \left(\mathbf{c}_{V1,dst}^{k,l} \oplus \mathbf{c}_{V1,src}^{k,l} \right), \quad (23)$$

$$\Delta_{V1}^{k,l} = \mathbf{C}'_{V1}{}^{k,l} \odot \frac{\partial \mathcal{L}}{\partial \mathbf{L}^{k,l}}. \quad (24)$$

This Δ matrix is never used as a matrix, but instead is summed along its rows and columns. Depending on which dimension it is summed over, it will correspond to gradient passed along source or destination edges. We can define the following vectors

$$\begin{aligned} \Sigma_{V1,dst}^{k,l} &= \text{Sum} \left(\Delta^{k,l}, \text{rows} \right), \\ \Sigma_{V1,src}^{k,l} &= \text{Sum} \left(\Delta^{k,l}, \text{cols} \right). \end{aligned} \quad (25)$$

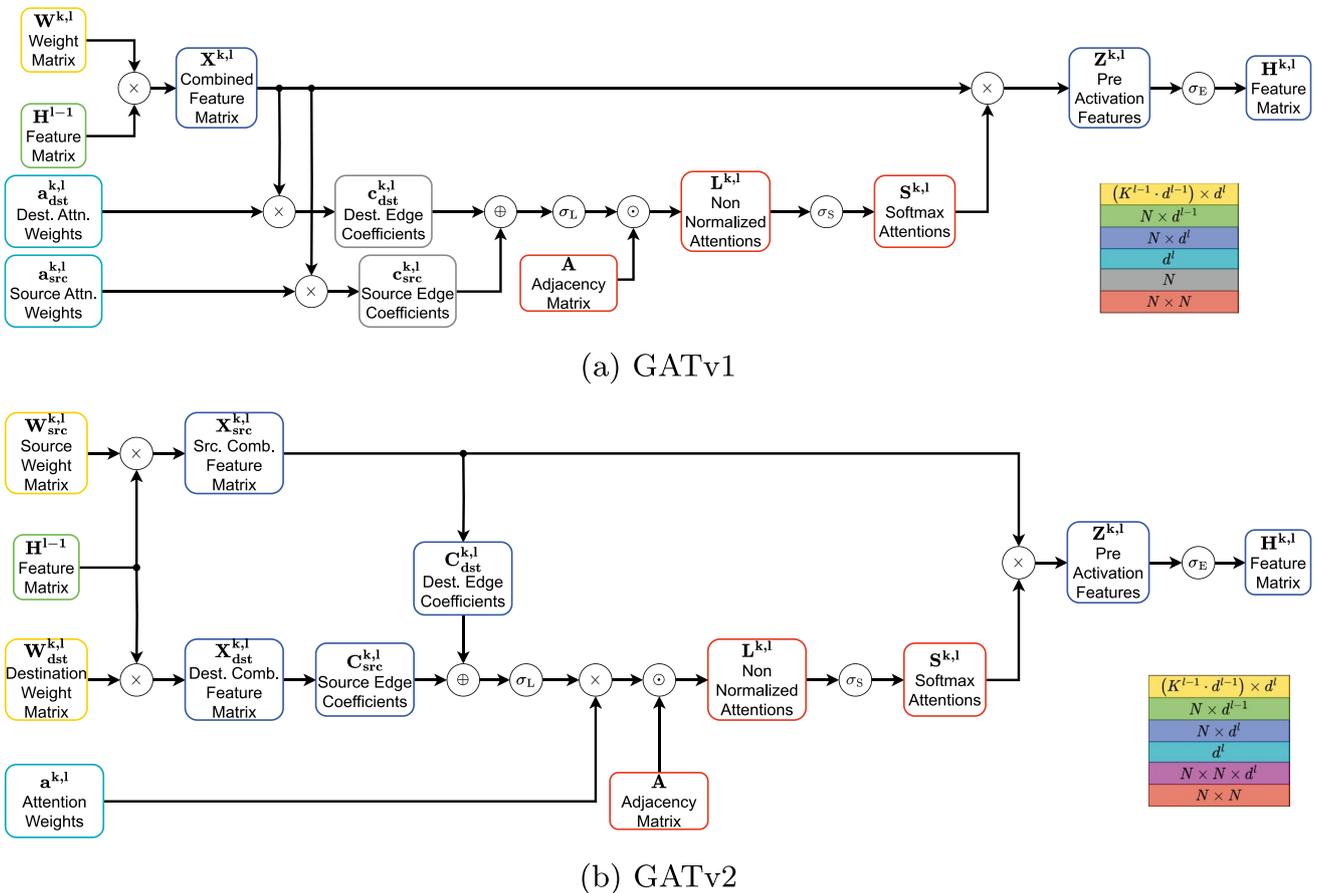


Figure 2 DFG for the forward pass of a single attention head within a hidden layer. Each box is colored according to the shape of the matrix at that point in the computation.

Table 2 Notations for GAT Training.

Notation		Description
$\sigma'_L, \sigma'_S, \sigma'_E$		Derivatives of LReLU, softmax, and ELU, respectively
\odot		Hadamard (element-wise) vector or matrix product
Notation	Size	Description
\mathcal{L}	$\mathbb{R}^{N \times d^l}$	Loss from output of GAT network
δ^l, δ'^l	$\mathbb{R}^{N \times (K^l \cdot d^l)}$	Gradient into layer l before and after being transformed by the activation derivative. Equivalent to $\frac{\partial \mathcal{L}}{\partial \mathbf{H}^l}$ and $\frac{\partial \mathcal{L}}{\partial \mathbf{Z}^l}$, respectively.
$\frac{\partial \mathcal{L}}{\partial \mathbf{L}^{k,l}}, \frac{\partial \mathcal{L}}{\partial \mathbf{S}^{k,l}}$	$\mathbb{R}^{N \times N}$	Loss with respect to pre- and post-softmax attention coefficients
$\mathbf{C}'_{V1, k,l}$	$\mathbb{R}^{N \times N}$	GATv1: Matrix holding LReLU derivative of the source/destination coefficient sums from the forward pass
$\mathbf{C}'_{V2, k,l}$	$\mathbb{R}^{N \times N \times d^l}$	GATv2: Matrix of vectors holding LReLU derivative of the source/destination coefficient sums from the forward pass
$\Delta_{V1, k,l}$	$\mathbb{R}^{N \times N}$	GATv1: Scalar loss with respect to edges
$\Delta_{V2, k,l}$	$\mathbb{R}^{N \times N \times d^l}$	GATv2: Vector loss with respect to edges
$\Sigma_{V1, dst, k,l}, \Sigma_{V1, src, k,l}$	\mathbb{R}^N	GATv1: Scalar losses with respect to nodes being destinations and sources of edges.
$\Sigma_{V2, dst, k,l}, \Sigma_{V2, src, k,l}$	$\mathbb{R}^{N \times d^l}$	GATv2: Vector losses with respect to nodes being destinations and sources of edges
$\frac{\partial \mathcal{L}}{\partial \mathbf{a}_{dst}^{k,l}}, \frac{\partial \mathcal{L}}{\partial \mathbf{a}_{src}^{k,l}}, \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{k,l}}$	\mathbb{R}^{d^l}	Loss with respect to the attention weights of head k of layer l
$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{k,l}}, \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{dst}^{k,l}}, \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{src}^{k,l}}$	$\mathbb{R}^{d^{l-1} \times d^l}$	Loss with respect to weights of head k of layer l

We can then derive the gradient of the loss with respect to the attention weights using these vectors

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{a}_{dst}^{k,l}} &= \mathbf{X}^{k,lT} \times \Sigma_{V1, dst}^{k,l} \\ \frac{\partial \mathcal{L}}{\partial \mathbf{a}_{src}^{k,l}} &= \mathbf{X}^{k,lT} \times \Sigma_{V1, src}^{k,l} \end{aligned} \tag{26}$$

These source and edge gradient vectors will also be used in computing the gradient of the loss with respect to input features and weights.

4.2.2 GATv2 Attention Gradient

The derivation for GATv2 is similar to GATv1, though here the $\mathbf{C}'^{k,l}$ matrix is instead a 3-dimensional tensor, because the attention coefficients are vectors instead of scalars. Effectively, it is a matrix identical in behavior to the GATv1 case, with the scalar elements replaced by vectors. This is written as

$$\mathbf{C}'_{V2, ij, k,l} = \mathbf{a}^{k,l} \odot \sigma'_L \left(\mathbf{c}_{i, V2, dst}^{k,l} + \mathbf{c}_{j, V2, src}^{k,l} \right), \tag{27}$$

$$\Delta_{V2, ij, k,l} = \frac{\partial \mathcal{L}}{\partial \mathbf{L}_{ij}^{k,l}} \cdot \mathbf{C}'_{V2, ij, k,l}, \tag{28}$$

$$\Delta_{V2, k,l} = \mathbf{C}'_{V2, k,l} \odot \frac{\partial \mathcal{L}}{\partial \mathbf{L}^{k,l}},$$

$$\begin{aligned} \Sigma_{V2, dst, k,l} &= \text{Sum}(\Delta_{V2, k,l}, \text{rows}), \\ \Sigma_{V2, src, k,l} &= \text{Sum}(\Delta_{V2, k,l}, \text{cols}). \end{aligned} \tag{29}$$

Instead of the $\Sigma^{k,l}$ tensors being vectors, they are matrices, owing to the vector elements in the $\Delta^{k,l}$ tensor. Finally, the loss with respect to attention weights can be derived as

$$\frac{\partial \mathcal{L}}{\partial \mathbf{a}^{k,l}} = \sum_{ij} \frac{\partial \mathcal{L}}{\partial \mathbf{L}_{ij}^{k,l}} \cdot \sigma_L \left(\mathbf{c}_{i, V2, dst}^{k,l} + \mathbf{c}_{j, V2, src}^{k,l} \right). \tag{30}$$

Unlike GATv1, this does not depend on the $\Sigma^{k,l}$ matrices, instead only depends on $\frac{\partial \mathcal{L}}{\partial \mathbf{L}^{k,l}}$. This is because the attention weights in GATv2 are not used until after the LReLU operation.

4.3 Gradients Associated with Combination

The final part of the backpropagation is similar between GATv1 and GATv2, consisting of combining the gradients for the feature and weight matrices with the contributions from the attention mechanism and the aggregation step. Because there are differences, we will again look at the two GAT versions separately.

4.3.1 GATv1 Weight and Feature Gradient

Both the feature and weight matrices can be written as a combination of gradient contributions from three sources: the gradients along the source and destination edges and the gradient due to weighted aggregation. The loss with respect to input features can be written as

$$\frac{\partial \mathcal{L}}{\partial \mathbf{H}^l} = \sum_{k=1}^{K^l} \left(\left(\boldsymbol{\Sigma}_{V1,dst}^{k,l} \times \mathbf{a}_{dst}^{k,lT} \times \mathbf{W}^{k,lT} \right) + \left(\boldsymbol{\Sigma}_{V1,src}^{k,l} \times \mathbf{a}_{src}^{k,lT} \times \mathbf{W}^{k,lT} \right) + \left(\mathbf{S}^{k,lT} \times \boldsymbol{\delta}'^{k,l} \times \mathbf{W}^{k,lT} \right) \right), \tag{31}$$

$$= \sum_{k=1}^{K^l} \left[\left(\boldsymbol{\Sigma}_{V1,dst}^{k,l} \times \mathbf{a}_{dst}^{k,lT} \right) + \left(\boldsymbol{\Sigma}_{V1,src}^{k,l} \times \mathbf{a}_{src}^{k,lT} \right) + \left(\mathbf{S}^{k,lT} \times \boldsymbol{\delta}'^{k,l} \right) \right] \times \mathbf{W}^{k,lT}, \tag{32}$$

$$= \boldsymbol{\delta}^{l-1}. \tag{33}$$

The loss with respect to weights can be written similarly

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{k,l}} = \left(\left(\mathbf{H}^{l-1T} \times \boldsymbol{\Sigma}_{V1,dst}^{k,l} \times \mathbf{a}_{dst}^{k,lT} \right) + \left(\mathbf{H}^{l-1T} \times \boldsymbol{\Sigma}_{V1,src}^{k,l} \times \mathbf{a}_{src}^{k,lT} \right) + \left(\mathbf{H}^{l-1T} \times \mathbf{S}^{k,lT} \times \boldsymbol{\delta}'^{k,l} \right) \right), \tag{34}$$

$$= \mathbf{H}^{l-1T} \times \left[\left(\boldsymbol{\Sigma}_{V1,dst}^{k,l} \times \mathbf{a}_{dst}^{k,lT} \right) + \left(\boldsymbol{\Sigma}_{V1,src}^{k,l} \times \mathbf{a}_{src}^{k,lT} \right) + \left(\mathbf{S}^{k,lT} \times \boldsymbol{\delta}'^{k,l} \right) \right]. \tag{35}$$

As can be seen from the equations, the terms in the square brackets are shared among both the gradients.

4.3.2 GATv2 Weight and Feature Gradient

The GATv2 equations are similar to GATv1, but because there are source and destination weight matrices, the three terms used to compute the feature gradient are split among the two weight matrix gradients. In GATv2, the loss with respect to input features is given by

$$\frac{\partial \mathcal{L}}{\partial \mathbf{H}^{l-1}} = \sum_{k=1}^{K^l} \left(\left(\boldsymbol{\Sigma}_{V2,dst}^{k,l} \times \mathbf{W}_{dst}^{k,lT} \right) + \left(\boldsymbol{\Sigma}_{V2,src}^{k,l} \times \mathbf{W}_{src}^{k,lT} \right) + \left(\mathbf{S}^{k,lT} \times \boldsymbol{\delta}'^{k,l} \times \mathbf{W}_{src}^{k,lT} \right) \right), \tag{36}$$

$$= \sum_{k=1}^{K^l} \left(\left(\boldsymbol{\Sigma}_{V2,dst}^{k,l} \times \mathbf{W}_{dst}^{k,lT} \right) + \left[\boldsymbol{\Sigma}_{V2,src}^{k,l} + \left(\mathbf{S}^{k,lT} \times \boldsymbol{\delta}'^{k,l} \right) \right] \times \mathbf{W}_{src}^{k,lT} \right), \tag{37}$$

$$= \boldsymbol{\delta}^{l-1}. \tag{38}$$

Then, the two weight matrix gradients are

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{dst}^{k,l}} = \mathbf{H}^{l-1T} \times \boldsymbol{\Sigma}_{V2,dst}^{k,l},$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{src}^{k,l}} = \left(\left(\mathbf{H}^{l-1T} \times \boldsymbol{\Sigma}_{V2,src}^{k,l} \right) + \left(\mathbf{H}^{l-1T} \times \mathbf{S}^{k,lT} \times \boldsymbol{\delta}'^{k,l} \right) \right),$$

$$= \mathbf{H}^{l-1T} \times \left[\boldsymbol{\Sigma}_{V2,src}^{k,l} + \left(\mathbf{S}^{k,lT} \times \boldsymbol{\delta}'^{k,l} \right) \right]. \tag{39}$$

The shared terms are once again shown in square brackets.

To summarize the backward pass pictorially, we modify the forward pass DFGs from Fig. 2 to show the dependencies during backpropagation, shown in Fig. 3. The example from the forward pass is used again to show the steps of backpropagation in Fig. 5.

5 Conclusion

This paper has presented closed-form matrix equations for the training of Graph Attention Networks. These equations explicitly illustrate the dependencies among various parts of the computations. These dependencies can be exploited to design accelerators for training GATs using techniques such as variable reuse, gradient interleaving and inter-layer pipelining [28].

Appendix. Supplementary Figures

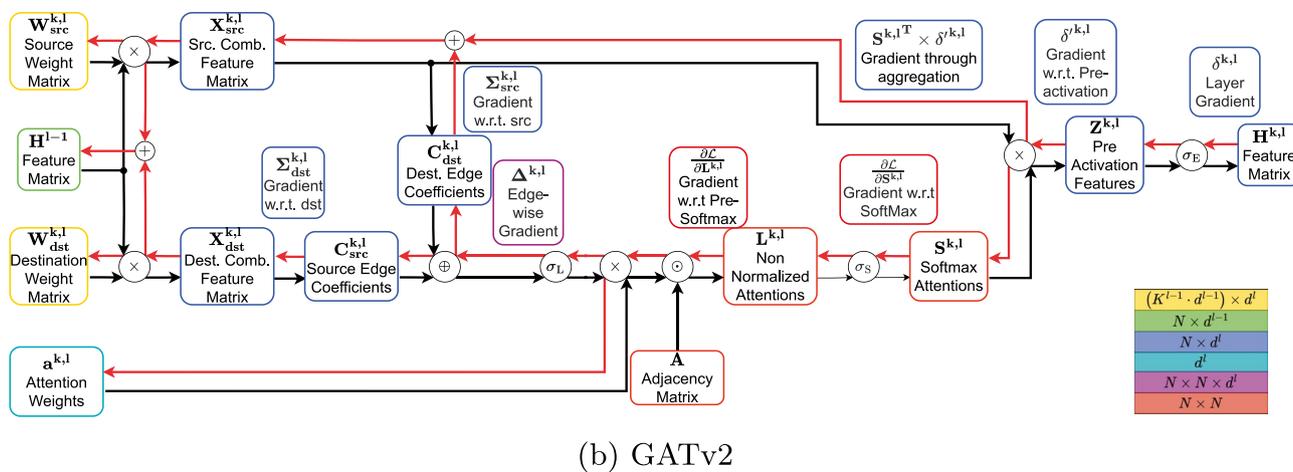
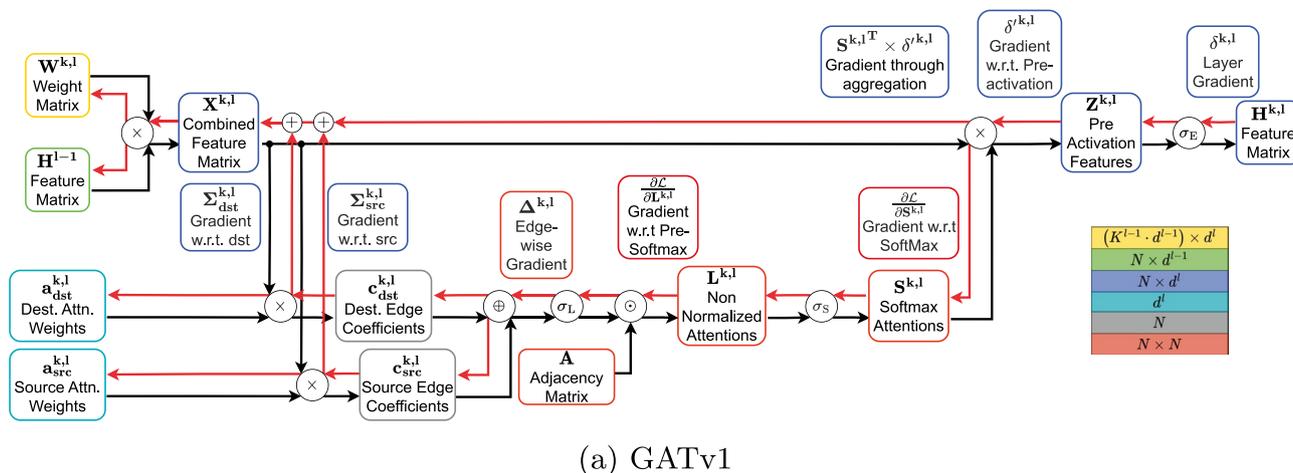


Figure 3 DFG for the forward and backward pass of a single attention head within a hidden layer. Backward pass matrices are shown adjacent to data flow. Colors of the boxes correspond to the shape of their respective matrices.

Figure 4 Forward pass of GAT Layer on an 8-Node graph, with 5-feature width input and 4-feature width attention head output. Differences between GATv1 and GATv2 are separated with a vertical line, such that GATv1’s operations appear on the left and GATv2’s on the right. The matrices are colored by shape according to the legends in Fig. 2. Subfigure **a** shows the input to the layer, with each node having a feature vector associated with it. Subfigure **b** shows the combination and edge coefficient calculation steps, corresponding to Eqs. (2)–(5), with each node having a source and destination coefficient associated with it. Subfigure **c** shows the pre-softmax attention coefficient calculation from Eqs. (8) and (9). Subfigure **d** shows the row-wise softmax operation on the matrix of these coefficients, from Eq. (10). Finally, subfigure **e** shows the aggregation step, weighted with the attention coefficients from Eqs. (11) and (12) (for brevity, we omit the near identical GATv2 equation using $\mathbf{X}_{src}^{k,l}$ instead of $\mathbf{X}^{k,l}$). The figure only shows a single attention head, and the output from subfigure **e** would be concatenated with the output of the other attention heads.

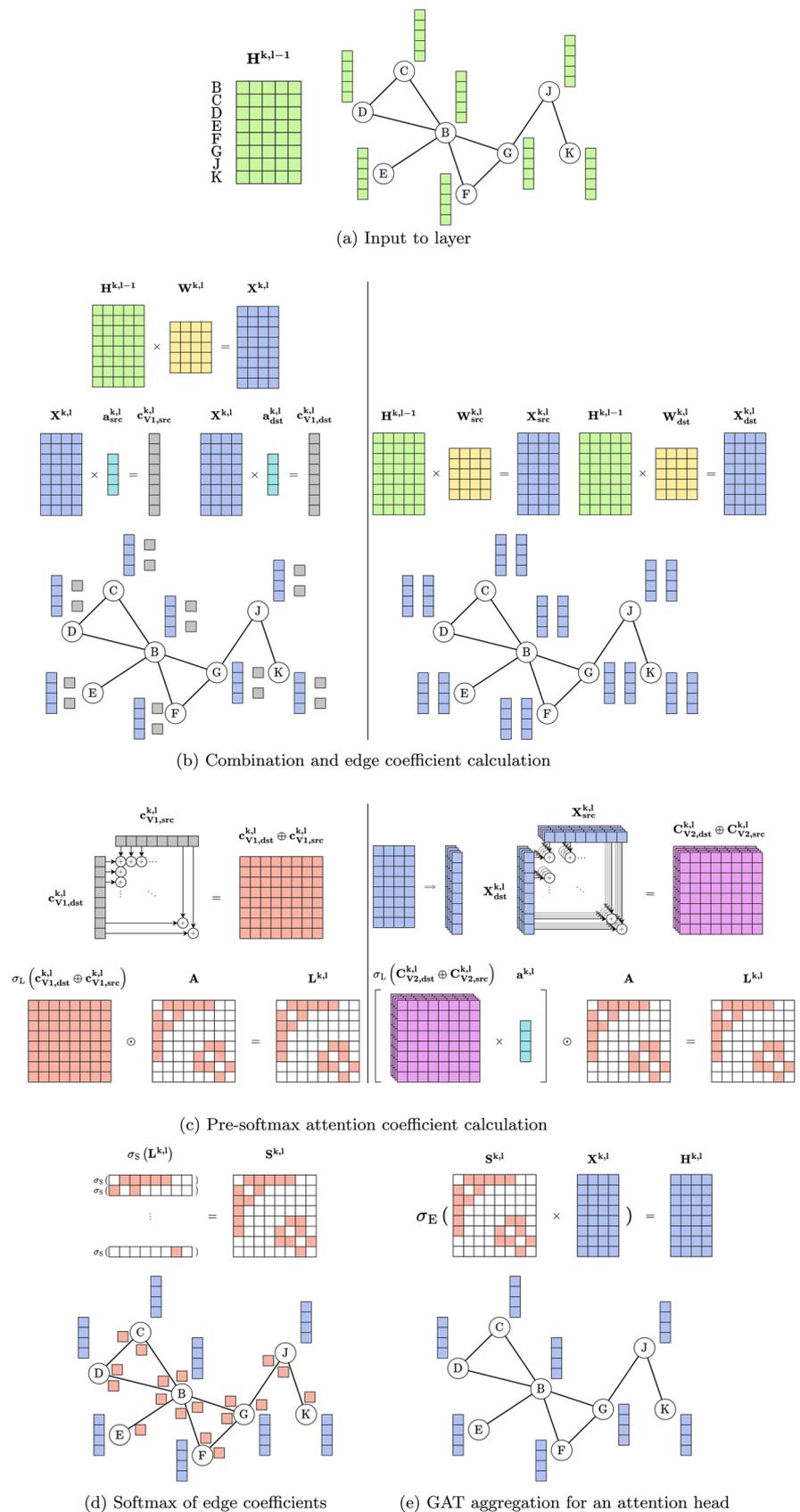


Figure 5 Backward pass for the forward pass example in Fig. 4. Differences between GATv1 and GATv2 are separated with a vertical line, such that GATv1’s operations appear on the left and GATv2’s on the right, except for the final step in weight and feature gradient calculation. The matrices are colored by shape according to the legends in Fig. 3. Subfigure **a** shows the calculation of the gradients with respect to the input of the layer activation function from Eq. (15). Subfigure **b** shows the gradient with respect to the aggregation step due to the attention coefficients, corresponding to Eq. (19). Subfigure **c** follows the gradient to the input of the softmax, covering Eq. (21). Subfigure **d** shows the computation of the gradient with respect to attention weights from Eqs. (26) and (30). In subfigure **e**, the $C^{k,l}$ matrix is shown following Eqs. (23) and (27) and used to find $\Delta^{k,l}$ as in Eqs. (24) and (28). Subfigure **f** shows the calculation of the $\Sigma^{k,l}$ matrices following Eqs. (25) and (29), and the multiplication with the attention weights in GATv1’s case. Finally, the gradient with respect to the weights and input features is shown for GATv1 in subfigure **g** from Eqs. (32) and (35) and GATv2 in subfigure **h** from Eqs. (32) and (39).

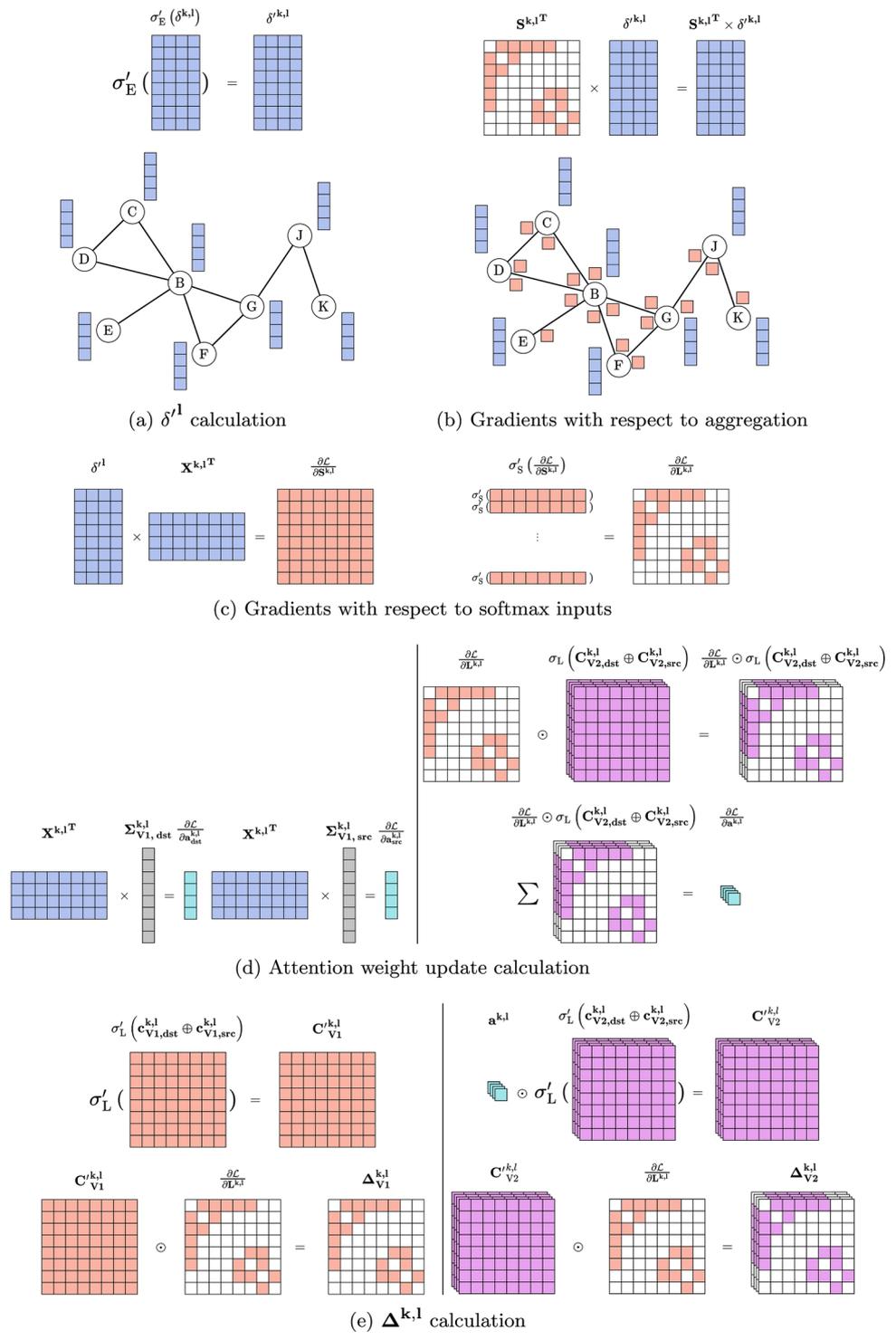
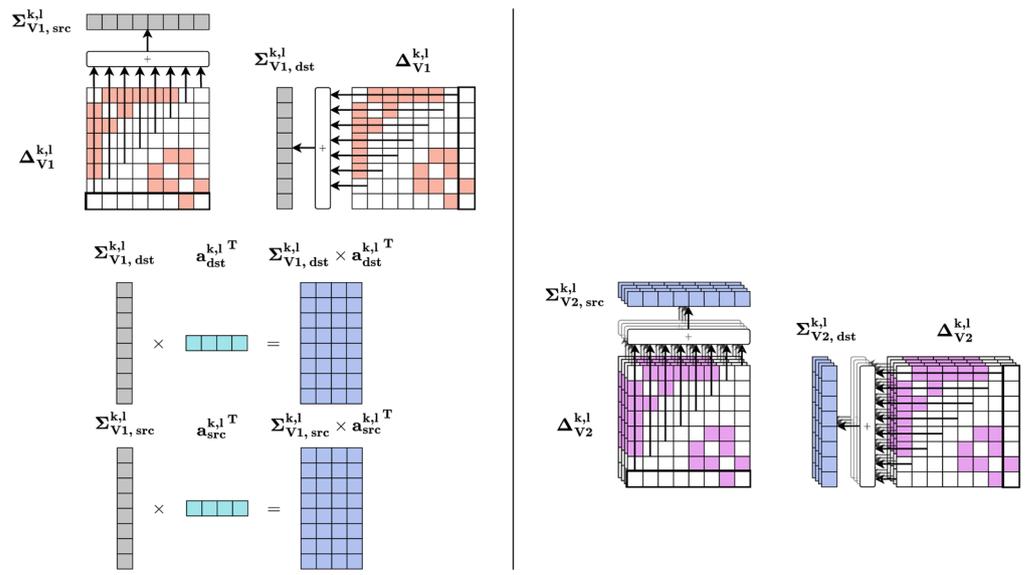


Figure 5 (continued)



(f) $\Sigma_{dst}^{k,l}$ and $\Sigma_{src}^{k,l}$ calculation

$$\sum_{k=1}^{K^l} \left[\Sigma_{V1,dst}^{k,l} \times a_{dst}^{k,l T} + \Sigma_{V1,src}^{k,l} \times a_{src}^{k,l T} \right] \times \mathbf{S}^{k,l T} \times \delta^{k,l} \times \mathbf{W}^{k,l T} = \delta^{l-1}$$

$$\mathbf{H}^{l-1 T} \times \left[\Sigma_{V1,dst}^{k,l} \times a_{dst}^{k,l T} + \Sigma_{V1,src}^{k,l} \times a_{src}^{k,l T} \right] \times \mathbf{S}^{k,l T} \times \delta^{k,l} = \frac{\partial \mathcal{L}}{\partial \mathbf{W}^{k,l T}}$$

(g) V1 weight and feature gradient calculation

$$\sum_{k=1}^{K^l} \left[\Sigma_{V2,dst}^{k,l} \times \mathbf{W}_{dst}^{k,l T} \right] + \left[\Sigma_{V2,src}^{k,l} + \mathbf{S}^{k,l T} \times \delta^{k,l} \right] \times \mathbf{W}_{src}^{k,l T} = \delta^{l-1}$$

$$\mathbf{H}^{l-1 T} \times \left[\mathbf{S}^{k,l T} \times \delta^{k,l} + \Sigma_{V2,src}^{k,l} \right] = \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{src}^{k,l T}} \quad \mathbf{H}^{l-1 T} \times \Sigma_{V2,dst}^{k,l} = \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{dst}^{k,l T}}$$

(h) V2 weight and feature gradient calculation

Acknowledgements The authors thank Nanda Unnikrishnan for numerous useful discussions.

Funding This paper was supported in part by the National Science Foundation under grant number CCF-1954749.

Data Availability Data sharing is not applicable to this article, as no datasets were generated or analyzed during the current study.

References

- Cheng, Z., Yan, C., Wu, F. X., & Wang, J. (2022). Drug-target interaction prediction using multi-head self-attention and graph attention network. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 19(4), 2208–2218. <https://doi.org/10.1109/TCBB.2021.3077905>. Conference Name: IEEE/ACM Transactions on Computational Biology and Bioinformatics.
- Yang, Z., Liu, J., Wang, Z., Wang, Y., & Feng, J. (2020). Multi-class metabolic pathway prediction by graph attention-based deep learning method. In: *2020 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pp. 126–131. <https://doi.org/10.1109/BIBM49941.2020.9313298>
- Zhao, H., Wang, Y., Duan, J., Huang, C., Cao, D., Tong, Y., Xu, B., Bai, J., Tong, J., & Zhang, Q. (2020). Multivariate time-series anomaly detection via graph attention network. In: *2020 IEEE International Conference on Data Mining (ICDM)*, pp. 841–850. <https://doi.org/10.1109/ICDM50108.2020.00093>. ISSN: 2374-8486.
- Zhang, C., James, J. Q., & Liu, Y. (2019). Spatial-temporal graph attention networks: A deep learning approach for traffic forecasting. *IEEE Access*, 7, 166246–166256. <https://doi.org/10.1109/ACCESS.2019.2953888>. Conference Name: IEEE Access.
- Balaji, S. S., & Parhi, K. K. (2023). Classifying Subjects with PFC Lesions from Healthy Controls during Working Memory Encoding via Graph Convolutional Networks. In: *2023 11th International IEEE/EMBS Conference on Neural Engineering (NER)*, pp. 1–4. <https://doi.org/10.1109/NER52421.2023.10123793>. ISSN: 1948-3554.
- Veličković, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., & Bengio, Y. (2018). Graph Attention Networks. *arXiv*. [arXiv:1710.10903](https://arxiv.org/abs/1710.10903) [cs, stat]. <http://arxiv.org/abs/1710.10903>. Accessed 24 Feb 2023.
- Brody, S., Alon, U., & Yahav, E. (January 2022). How attentive are graph attention networks? Technical Report [arXiv:2105.14491](https://arxiv.org/abs/2105.14491), [arXiv:2105.14491](https://arxiv.org/abs/2105.14491) [cs] type: article. <http://arxiv.org/abs/2105.14491>. Accessed 2024 Feb 2023.
- Unnikrishnan, N. K., & Parhi, K. K. (2023). InterGrad: Energy-efficient training of convolutional neural networks via interleaved gradient scheduling. *IEEE Transactions on Circuits and Systems I: Regular Papers*, 70(5), 1949–1962. <https://doi.org/10.1109/TCSI.2023.3246468>. Conference Name: IEEE Transactions on Circuits and Systems I: Regular Papers
- Gori, M., Monfardini, G., & Scarselli, F. (2005). A new model for learning in graph domains. In: *Proceedings. 2005 IEEE International Joint Conference on Neural Networks*, vol. 2, pp. 729–7342. <https://doi.org/10.1109/IJCNN.2005.1555942>. ISSN: 2161-4407
- Scarselli, F., Gori, M., Tsoi, A. C., Hagenbuchner, M., & Monfardini, G. (2009). The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1), 61–80. <https://doi.org/10.1109/TNN.2008.2005605>. Conference Name: IEEE Transactions on Neural Networks
- Wu, Z., Pan, S., Chen, F., Long, G., Zhang, C., & Yu, P. S. (2021). A comprehensive survey on graph neural networks. *IEEE Transactions Neural Networks Learning System*, 32(1), 4–24. <https://doi.org/10.1109/TNNLS.2020.2978386>. [arXiv:1901.00596](https://arxiv.org/abs/1901.00596) [cs, stat]. Accessed 24 Feb 2023.
- Parhi, K. K., & Unnikrishnan, N. K. (2020). Brain-inspired computing: models and architectures. *IEEE Open Journal of Circuits and Systems*, 1, 185–204. <https://doi.org/10.1109/OJCS.2020.3032092>. Conference Name: IEEE Open Journal of Circuits and Systems
- Kipf, T. N., & Welling, M. (February 2017). Semi-supervised classification with graph convolutional networks. *Technical Report*. [arXiv:1609.02907](https://arxiv.org/abs/1609.02907), [arXiv:1609.02907](https://arxiv.org/abs/1609.02907) [cs, stat] type: article. <http://arxiv.org/abs/1609.02907>. Accessed 24 Feb 2023.
- Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., & Dahl, G. E. (June 2017). Neural message passing for quantum chemistry. *Technical Report*. [arXiv:1704.01212](https://arxiv.org/abs/1704.01212), [arXiv:1704.01212](https://arxiv.org/abs/1704.01212) [cs] type: article. <http://arxiv.org/abs/1704.01212>. Accessed 24 Feb 2023.
- Zhang, B., & Prasanna, V. (2023). Dynaspase: Accelerating GNN inference through dynamic sparsity exploitation. <https://arxiv.org/abs/2303.12901v1>. Accessed 3 Jun 2023.
- Mondal, S., Manasi, S. D., Kunal, K., Ramprasath, S., & Sapatnekar, S. S. (2021). GNNIE: GNN inference engine with load-balancing and graph-specific caching. <https://arxiv.org/abs/2105.10554v2>. Accessed 3 Jun 2023.
- He, Z., Tian, T., Wu, Q., & Jin, X. (2023). FTW-GAT: An FPGA-based accelerator for graph attention networks with ternary weights. *IEEE Transactions on Circuits and Systems II: Express Briefs*, 1–1. <https://doi.org/10.1109/TCSII.2023.3280180>. Conference Name: IEEE Transactions on Circuits and Systems II: Express Briefs.
- Geng, T., Wu, C., Zhang, Y., Tan, C., Xie, C., You, H., Herbordt, M., Lin, Y., & Li, A. (2021). I-GCN: A graph convolutional network accelerator with runtime locality enhancement through islandization. In: *MICRO-54: 54th Annual IEEE/ACM International Symposium on Microarchitecture*. *MICRO '21*, pp. 1051–1063. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3466752.3480113>
- Zeng, H., & Prasanna, V. (2019). GraphACT: Accelerating GCN training on CPU-FPGA heterogeneous platforms. <https://doi.org/10.1145/3373087.3375312>. <https://arxiv.org/abs/2001.02498v1>. Accessed 3 Jun 2023.
- Chen, X., Wang, Y., Xie, X., Hu, X., Basak, A., Liang, L., Yan, M., Deng, L., Ding, Y., Du, Z., & Xie, Y. (2022). Rubik: A hierarchical architecture for efficient graph neural network training. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 41(4), 936–949. <https://doi.org/10.1109/TCAD.2021.3079142>. Conference Name: IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems.
- Zheng, D., Ma, C., Wang, M., Zhou, J., Su, Q., Song, X., Gan, Q., Zhang, Z., & Karypis, G. (2020). DistDGL: Distributed graph neural network training for billion-scale graphs. <https://arxiv.org/abs/2010.05337v3>. Accessed 3 Jun 2023.
- Lin, Z., Li, C., Miao, Y., Liu, Y., & Xu, Y. (2020). Pagraph: Scaling GNN training on large graphs via computation-aware caching. In: *Proceedings of the 11th ACM Symposium on Cloud Computing*. *SoCC '20*, pp. 401–415. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3419111.3421281>
- Lin, Y.-C., Zhang, B., & Prasanna, V. (2023). HitGNN: High-throughput GNN training framework on CPU+Multi-FPGA heterogeneous platform. <https://arxiv.org/abs/2303.01568v1>. Accessed 3 Jun 2023.
- Luong, M.-T., Pham, H., Manning, C.D. (2015). Effective approaches to attention-based neural machine translation. [arXiv:1508.04025](https://arxiv.org/abs/1508.04025) [cs]. <http://arxiv.org/abs/1508.04025>. Accessed 27 Jan 2023.
- Gehring, J., Auli, M., Grangier, D., & Dauphin, Y. N. (July 2017). A convolutional encoder model for neural machine translation. *Technical Report*. [arXiv:1611.02344](https://arxiv.org/abs/1611.02344), [arXiv:1611.02344](https://arxiv.org/abs/1611.02344)

- [cs] type: article. <http://arxiv.org/abs/1611.02344>. Accessed 27 Jan 2023.
26. Mnih, V., Heess, N., Graves, A., & Kavukcuoglu, K. (June 2014). Recurrent Models of Visual Attention. *Technical Report*. [arXiv:1406.6247](https://arxiv.org/abs/1406.6247), arXiv. [arXiv:1406.6247](https://arxiv.org/abs/1406.6247) [cs, stat] type: article. <http://arxiv.org/abs/1406.6247>. Accessed 27 Jan 2023.
 27. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (December 2017). Attention is all you need. *Technical Report*. [arXiv:1706.03762](https://arxiv.org/abs/1706.03762), arXiv. [arXiv:1706.03762](https://arxiv.org/abs/1706.03762) [cs] type: article. <http://arxiv.org/abs/1706.03762>. Accessed 27 Jan 2023.
 28. Unnikrishnan, N. K., & Parhi, K. K. (2021). LayerPipe: Accelerating deep neural network training by intra-layer and inter-layer gradient pipelining and multiprocessor scheduling. In: *2021 IEEE/ACM International Conference On Computer Aided Design (ICCAD)*, pp. 1–8. <https://doi.org/10.1109/ICCAD51958.2021.9643567>. ISSN: 1558-2434.

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.



Joe Gould (Student member, IEEE) is currently pursuing an M.S. degree in electrical engineering at the University of Minnesota, Minneapolis, USA. He worked at Commscope, USA, from 2020 to 2021 under a co-op internship program as a design verification engineer for FPGA products. His interests are in the design of efficient architectures for accelerator and signal processing systems.



Keshab K. Parhi (Fellow, IEEE) received the B.Tech. degree from the Indian Institute of Technology (IIT), Kharagpur, in 1982, the M.S.E.E. degree from the University of Pennsylvania, Philadelphia, in 1984, and the Ph.D. degree from the University of California, Berkeley, in 1988. He has been with the University of Minnesota, Minneapolis, since 1988, where he is currently Erwin A. Kelen Chair and Distinguished McKnight University Professor in the Department of Electrical and Computer Engineering. He has published over 700 papers, is the inventor of 34 patents,

and has authored the textbook *VLSI Digital Signal Processing Systems* (Wiley, 1999). His current research addresses VLSI architecture design of machine learning and signal processing systems, hardware security, and data-driven neuroengineering and neuroscience. Dr. Parhi is the recipient of numerous awards including the 2017 Mac Van Valkenburg award and the 2012 Charles A. Desoer Technical Achievement award from the IEEE Circuits and Systems Society, the 2003 IEEE Kiyo Tomiyasu Technical Field Award, and a Golden Jubilee medal from the IEEE Circuits and Systems Society in 2000. He served as the Editor-in-Chief of the *IEEE Trans. Circuits and Systems, Part-I* during 2004 and 2005. He is a Fellow of the American Association for the Advancement of Science (AAAS), the Association for Computing Machinery (ACM), the American Institute of Medical and Biological Engineering (AIMBE), and the National Academy of Inventors (NAI).