

# GENERALIZED ATTENTION FLOW: FEATURE ATTRIBUTION FOR TRANSFORMER MODELS VIA MAXIMUM FLOW

Anonymous authors  
Paper under double-blind review

## ABSTRACT

This paper introduces Generalized Attention Flow, a novel feature attribution method for Transformer models that addresses the limitations of existing approaches. By generalizing Attention Flow and substituting attention weights with an arbitrary Information Tensor, the method leverages attention weights, their gradients, maximum flow, and the barrier method to generate more accurate feature attributions. The proposed approach demonstrates important theoretical properties and resolves issues associated with previous methods that rely solely on simple aggregation of attention weights. Comprehensive benchmarking in NLP sequence classification tasks reveals that a specific variant of Generalized Attention Flow consistently outperforms state-of-the-art feature attribution methods across most evaluation scenarios, offering a more accurate explanation of Transformer model outputs.

## 1 INTRODUCTION

Feature attribution methods are essential for building interpretable machine learning models. These methods assign a score to each input feature, reflecting its importance to the model’s output, thereby facilitating the understanding of model predictions.

The rise of Transformer models with self-attention mechanisms has necessitated feature attribution methods for interpreting these models (Vaswani et al., 2017; Bahdanau et al., 2016; Devlin et al., 2019; Sanh et al., 2020; Dosovitskiy et al., 2020; Kobayashi et al., 2021). Initially, attention weights were considered potential feature attributions, but recent studies have questioned their effectiveness in explaining deep neural networks (Abnar & Zuidema, 2020; Clark et al., 2019; Jain & Wallace, 2019; Serrano & Smith, 2019). Consequently, various post hoc methods have been developed to obtain feature attributions for Transformer models.

Recent advancements in XAI have introduced numerous gradient-based methods, including Grads and AttGrads (Barkan et al., 2021), which leverage saliency to interpret Transformer outputs. Qiang et al. (2022) proposed AttCAT, integrating features, their gradients, and attention weights to quantify input influence on model outputs. However, many of these techniques still focus primarily on gradients of attention weights, inheriting limitations of earlier attention-based approaches.

Layer-wise Relevance Propagation (LRP) (Bach et al., 2015; Voita et al., 2019) transfers relevance scores from output to input. Chefer et al. (2021a;b) proposed a comprehensive method enabling information propagation through all Transformer components. However, this approach relies on specific LRP rules, limiting its applicability across various Transformer architectures.

Many existing methods for evaluating feature attributions in Transformers fail to capture pairwise interactions among features. This limitation arises from the independent computation of importance scores, which neglects feature interactions. For example, when calculating gradients of attention weights, they propagate directly from the output to the individual input feature, ignoring interactions. Additionally, many methods used to compute feature attributions in Transformers violate key axioms such as symmetry, sensitivity, efficiency, and linearity (Shapley, 1952; Sundararajan et al., 2017; Sundararajan & Najmi, 2020) (Sec. 3.5).

047 Recently, [Abnar & Zuidema \(2020\)](#) introduced Attention Flow to overcome these limitations in XAI methods.  
048 Attention Flow considers attention as capacities in a maximum flow problem, determining feature attributions  
049 based on the solution. This approach naturally captures the influence of attention mechanisms, as the  
050 paths of high attention through a network correspond to the flow of information from features to outputs.  
051 Applicable to any encoder-only Transformer, Attention Flow has demonstrated strong potential to improve  
052 model interpretability ([Abnar & Zuidema, 2020](#); [Modarressi et al., 2023](#); [Kobayashi et al., 2021](#)).

053 Subsequently, [Ethayarajh & Jurafsky \(2021\)](#) sought to connect attention flows and XAI by utilizing Shapley  
054 values ([Shapley, 1952](#)). While they aimed to show that attention flows could be interpreted as Shapley values  
055 under certain conditions, they overlooked the issue of non-uniqueness in such flows ([Sec. 3.3](#)).

056 **Our contributions.** We propose Generalized Attention Flow, which satisfies important theoretical properties  
057 and enhanced empirical performance. Specifically, our contributions are:

- 058 1. We introduce Generalized Attention Flow, an extension of the previously described Attention Flow. In this  
059 approach, feature attributions are generated by using the log barrier method to solve a regularized maximum  
060 flow problem within a capacity network formed from the functions applied to attention weights. Instead of  
061 defining capacities based solely on attention weights, we suggest using the gradients of these weights (GF) or  
062 the product of attention weights and their gradients (AGF) as alternatives.
- 063 2. We have addressed the non-uniqueness issue in Attention Flow, which invalidates some of its previously  
064 suggested theoretical properties ([Ethayarajh & Jurafsky, 2021](#)). Furthermore, we show that non-unique  
065 solutions occur frequently in practice. We have introduced barrier regularization to mitigate this issue to  
066 ensure a unique solution. As a result, we have demonstrated that feature attributions derived from the  
067 regularized maximum flow problem align with Shapley values and satisfy the axioms of efficiency, symmetry,  
068 nullity, and linearity ([Shapley, 1952](#); [Young, 1985](#); [Chen et al., 2023b](#)).
- 069 3. We extensively benchmarked the proposed feature attribution methods, defined using Generalized Attention  
070 Flow, against various existing state-of-the-art attribution methods. We found that a type of the proposed  
071 attribution methods outperforms previous state-of-the-art methods in terms of explanation performance for  
072 classification tasks across most evaluation scenarios, as measured by AOPC ([Barkan et al., 2021](#); [Nguyen,](#)  
073 [2018](#); [Chen et al., 2020](#)), LOdds ([Chen et al., 2020](#); [Shrikumar et al., 2018](#)), and classification metrics.
- 074 4. We have developed an open-source Python package for calculating feature attributions using Generalized  
075 Attention Flow. This package is distinctively flexible, capable of being applied to any encoder-only Trans-  
076 former model available in the [Hugging Face Transformers](#) package ([Wolf et al., 2020](#)). Furthermore, our  
077 methods are easily adaptable for various NLP tasks.

## 078 2 PRELIMINARIES

### 079 2.1 MULTI-HEAD ATTENTION MECHANISM

080 Given the input sequence  $\mathbf{X} \in \mathbb{R}^{t \times d}$ , where  $d$  is the dimensionality of the model’s input vectors and  $t$  is the  
081 number of tokens, the multi-head self-attention mechanism computes attention weights for each element in  
082 the sequence employing the following steps:

- 083 • **Linear Transformation:**

$$084 \mathbf{Q}_i = \mathbf{XW}_i^Q, \quad \mathbf{K}_i = \mathbf{XW}_i^K, \quad \mathbf{V}_i = \mathbf{XW}_i^V \quad (1)$$

085 Here  $\mathbf{Q}_i, \mathbf{K}_i \in \mathbb{R}^{t \times d_k}$  and  $\mathbf{V}_i \in \mathbb{R}^{t \times d_v}$ , where  $d_k$  and  $d_v$  represent the dimensionality of the key  
086 vector and value vector respectively, and  $i$  represents the index of the attention head.

094 • **Scaled Dot-Product Attention:**

095 
$$\text{Attention}_i(\mathbf{Q}_i, \mathbf{K}_i, \mathbf{V}_i) = \tilde{\mathbf{A}}_i \mathbf{V}_i \quad (2)$$

096 where the matrix of attention weights  $\tilde{\mathbf{A}}_i \in \mathbb{R}^{t \times t}$  is defined as:

097 
$$\tilde{\mathbf{A}}_i = \text{softmax} \left( \frac{\mathbf{Q}_i \mathbf{K}_i^T}{\sqrt{d_k}} \right) \quad (3)$$

101 • **Concatenation and Linear Projection:**

102 
$$\text{MultiHead}(\mathbf{X}) = \text{Concat}(\text{Attention}_1, \text{Attention}_2, \dots, \text{Attention}_h) \mathbf{W}^O \quad (4)$$

103 where  $\text{MultiHead}(\mathbf{X}) \in \mathbb{R}^{t \times d}$  and  $\mathbf{W}^O \in \mathbb{R}^{h \cdot d_v \times d}$ .

104 For a Transformer with  $l$  attention layers, the attention weights at each layer can be defined as multi-head attention weights:

105 
$$\hat{\mathbf{A}} = \text{Concat}(\tilde{\mathbf{A}}_1, \tilde{\mathbf{A}}_2, \dots, \tilde{\mathbf{A}}_h) \in \mathbb{R}^{h \times t \times t} \quad (5)$$

106 Extending this to a Transformer architecture itself, the Transformer attention weights  $\mathbf{A}$  can be defined as:

107 
$$\mathbf{A} = \text{Concat}(\hat{\mathbf{A}}_1, \hat{\mathbf{A}}_2, \dots, \hat{\mathbf{A}}_l) \in \mathbb{R}^{l \times h \times t \times t} \quad (6)$$

108 where  $\hat{\mathbf{A}}_j \in \mathbb{R}^{h \times t \times t}$  is the multi-head attention weight for the  $j$ -th attention layer.

112 2.2 MINIMUM-COST CIRCULATION & MAXIMUM FLOW PROBLEM

113 **Definition 2.1 (Minimum Cost Circulation).** Given a network  $G = (V, E, \mathbf{u}, \mathbf{l}, \mathbf{c})$  with  $|V| = n$  vertices and  $|E| = m$  edges, where  $c_{ij}$  is the cost,  $l_{i,j}$  and  $u_{i,j}$  are respectively the lower and upper capacities (demands) for the edge  $(i, j) \in E$ , circulation is a function  $f : E \rightarrow \mathbb{R}^{\geq 0}$  s.t.

114 
$$\begin{aligned} l_{ij} \leq f_{ij} \leq u_{ij}, & \quad \forall (i, j) \in E \\ \sum_{j:(i,j) \in E} f_{ij} - \sum_{j:(j,i) \in E} f_{ji} = 0, & \quad \forall i \in V. \end{aligned} \quad (7)$$

115 The min-cost circulation problem is to find a circulation  $f$  minimizing the cost function  $\sum_{(i,j) \in E} c_{ij} f_{ij}$ .

116 The minimum-cost circulation problem can be algebraically written as the following primal-dual linear programming (LP) problem (Van Den Brand et al., 2021; Chen et al., 2023a):

117 (Primal)  $\arg \min_{\substack{\mathbf{B}^\top \mathbf{f} = \mathbf{0} \\ l_e \leq f_e \leq u_e \forall e \in E}} \mathbf{c}^\top \mathbf{f}$  i.e.  $\arg \min_{\substack{\mathbf{B}^\top \mathbf{f} = \mathbf{0} \\ l \leq \mathbf{f} \leq \mathbf{u}}} \mathbf{c}^\top \mathbf{f}$ , (Dual)  $\arg \max_{\mathbf{B}\mathbf{y} + \mathbf{s} = \mathbf{c}} \sum_i \min(l_i s_i, u_i s_i)$  (8)

118 where  $\mathbf{B}_{m \times n}$ , is the edge-vertex incidence matrix. For a directed graph, the entries of the matrix  $\mathbf{B}$  are defined as follows:

119 
$$\mathbf{B}_{ev} = \begin{cases} -1, & \text{if vertex } v \text{ is the tail of edge } e, \\ 1, & \text{if vertex } v \text{ is the head of edge } e, \\ 0, & \text{if edge } v \text{ is not incident to vertex } e. \end{cases}$$

120 **Remark 2.1.** The maximum flow problem can be considered as a specific minimum-cost circulation problem. Here,  $\mathbf{B}$  is an edge-vertex incidence matrix of the input graph after we added to it an edge  $e(t, s)$  that connects the target  $t$  to the source  $s$  and its lower capacity  $l_{t,s}$  be 0 and its upper capacity  $u_{t,s}$  be  $\|\mathbf{u}\|_1$ . Also, the cost vector  $\mathbf{c}$  is a vector in which  $c_{t,s} = -1$  and  $c_e = 0$  for all other edges  $e \in E$  (Cormen et al., 2009).

141 2.3 BARRIER METHODS FOR CONSTRAINED OPTIMIZATION

142 Consider the following optimization problem:

143 
$$f^* = \arg \min_{\substack{\alpha(f)=0 \\ \beta(f) \leq 0}} \xi(f) \tag{9}$$

144 where  $h$  represents a convex inequality constraint,  $g$  represents an affine equality constraint, and  $f^*$  denote

145 the optimal solution.

146 The interior of the constraint region is defined as  $S = \{f \mid \alpha(f) = 0, \beta(f) < 0\}$ . Assuming  $S$  is nonempty

147 and convex, we introduce a barrier function  $\psi(f)$  on  $S$  that is continuous and approaches infinity as  $f$

148 approaches to the boundary of the region, specifically  $\lim_{\beta(f) \rightarrow 0^-} \psi(f) = \infty$ . One common example of

149 barrier functions is the log barrier function, which is represented as  $\log(-\beta(f))$ .

150 Given a barrier function  $\psi(f)$ , we can define a new objective function  $\xi(f) + \mu\psi(f)$ , where  $\mu$  is a positive

151 real number, which enables us to eliminate the inequality constraints in the original problem and obtain the

152 following problem:

153 
$$f_\mu^* = \arg \min_{\alpha(f)=0} \xi(f) + \mu\psi(f) \tag{10}$$

154 **Theorem 2.1.** *For any strictly convex barrier function  $\psi(f)$ , convex function  $\xi(f)$ , and  $\mu > 0$ , there exists a*

155 *unique optimal point  $f_\mu^*$ . Furthermore,  $\lim_{\mu \rightarrow 0} f_\mu^* = f^*$ , indicating that for any arbitrary  $\epsilon > 0$ , we can*

156 *select a sufficiently small  $\mu > 0$  such that  $\|f_\mu^* - f^*\| < \epsilon$  (van den Brand et al., 2023).*

157

158

159

160

161

162

163 3 METHODS

164

165 3.1 INFORMATION TENSOR

166 In Transformer-based networks, information propagation occurs through pathways facilitated by the attention

167 mechanism. These pathways can be conceptualized as routes within a graph structure, where tokens are

168 represented by nodes and computations are denoted by edges. The capacities of these edges correspond to

169 meaningful computational quantities that reflect the flow of information through the network (Ferrando &

170 Voita, 2024; Mueller, 2024).

171 Attention scores can represent the flow of information through the neural network during the feed-forward

172 phase of training, quantifying the importance of different input parts in generating the output (Abnar &

173 Zuidema, 2020; Ferrando & Voita, 2024). Additionally, the gradient of attention scores captures the flow of

174 information during back-propagation, reflecting how changes in the output influence the attention mechanism

175 throughout the network (Barkan et al., 2021). A combined view of attention scores and their gradients can

176 simultaneously represent information circulation during both feed-forward and back-propagation, offering a

177 comprehensive perspective on the network’s information dynamics (Barkan et al., 2021; Qiang et al., 2022;

178 Chefer et al., 2021a;b).

179 Our Generalized Attention Flow builds on this foundation, using an information tensor  $\bar{A} \in \mathbb{R}^{l \times t \times t}$  to

180 aggregate Transformer attention weights  $A$ , as defined in eq. 6. Based on the insights above, we propose

181 three aggregation functions to generate information tensors (Barkan et al., 2021; Chefer et al., 2021a):

- 182
- 183 1. **Attention Flow (AF):**  $\bar{A} := \mathbb{E}_h(A)$
  - 184 2. **Attention Grad Flow (GF):**  $\bar{A} := \mathbb{E}_h([\nabla A]_+)$
  - 185 3. **Attention  $\times$  Attention Grad Flow (AGF):**  $\bar{A} := \mathbb{E}_h([A \odot \nabla A]_+)$
- 186

187 Here,  $[x]_+ = \max(x, 0)$ ,  $\odot$  represents the Hadamard product,  $\nabla A := \frac{\partial y_t}{\partial A}$  where  $y_t$  is the model’s scalar

output, and  $\mathbb{E}_h$  denotes the mean across attention heads.

188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234

---

### Algorithm 1 Backward Information Capacity

---

**Input:**  $\bar{\mathbf{A}}_{l \times t \times t}$ : An information tensor.  
**Output:** Tuple:  $(\mathcal{A}, \mathbf{l}, \tilde{\mathbf{l}}, \mathbf{u}, \tilde{\mathbf{u}}, ss, st)$   
**function** GET\_BACKWARD\_CAPACITY( $\bar{\mathbf{A}}$ )  
 ▷ Initialization  
 $l, t, _ \leftarrow \bar{\mathbf{A}}.shape()$   
 $\beta_{\min} \leftarrow \min(\bar{\mathbf{A}} > 0)$   
 $\beta \leftarrow -\lfloor \log_{10}(\beta_{\min}) \rfloor$   
 $\gamma \leftarrow 10^\beta$   
 $Q_{tl} \leftarrow t * (l + 1) + 2$   
 $\mathbf{l} \leftarrow \text{zeros}(Q_{tl}, Q_{tl})$   
 $\mathbf{u} \leftarrow \text{zeros}(Q_{tl}, Q_{tl})$   
 $u_\infty \leftarrow t$   
 ▷ Fill super-source  $\rightarrow$  First Layer  
**for**  $i$  in range( $t$ ) **do**  
    $\mathbf{u}[i + 1][0] \leftarrow u_\infty$   
**end for**  
 ▷ Fill Last Layer  $\rightarrow$  super-target  
**for**  $i$  in range( $t$ ) **do**  
    $\mathbf{u}[-1][-i - 2] \leftarrow u_\infty$   
**end for**  
 ▷ Fill  $j$ -th Layer to  $(j + 1)$ -th Layer  
**for**  $j$  in range( $l$ ) **do**  
    $start \leftarrow t * j + 1$   
    $mid \leftarrow t * (j + 1) + 1$   
    $end \leftarrow t * (j + 2) + 1$   
    $\mathbf{u}[mid:end, start:mid] \leftarrow \bar{\mathbf{A}}_{[j,:,:,]}$   
**end for**  
 ▷ Get Integral Version of Capacities  
 $\tilde{\mathbf{l}} \leftarrow \text{int}(\gamma * \mathbf{l})$   
 $\tilde{\mathbf{u}} \leftarrow \text{int}(\gamma * \mathbf{u})$   
 ▷ Get Adjacency Matrix  
 $\mathcal{A} \leftarrow \mathbb{I}_{(\mathbf{u} > 0)}$   
 ▷ Get super-source and super-target  
 $ss, st \leftarrow t * (l + 1) + 1, 0$   
**end function**

---



---

### Algorithm 2 Forward Information Capacity

---

**Input:**  $\bar{\mathbf{A}}_{l \times t \times t}$ : An information tensor.  
**Output:** Tuple:  $(\mathcal{A}, \mathbf{l}, \tilde{\mathbf{l}}, \mathbf{u}, \tilde{\mathbf{u}}, ss, st)$   
**function** GET\_FORWARD\_CAPACITY( $\bar{\mathbf{A}}$ )  
 ▷ Initialization  
 $l, t, _ \leftarrow \bar{\mathbf{A}}.shape()$   
 $\beta_{\min} \leftarrow \min(\bar{\mathbf{A}} > 0)$   
 $\beta \leftarrow -\lfloor \log_{10}(\beta_{\min}) \rfloor$   
 $\gamma \leftarrow 10^\beta$   
 $Q_{tl} \leftarrow t * (l + 1) + 2$   
 $\mathbf{l} \leftarrow \text{zeros}(Q_{tl}, Q_{tl})$   
 $\mathbf{u} \leftarrow \text{zeros}(Q_{tl}, Q_{tl})$   
 $u_\infty \leftarrow t$   
 ▷ Fill super-source  $\rightarrow$  First Layer  
**for**  $i$  in range( $t$ ) **do**  
    $\mathbf{u}[0][i + 1] \leftarrow u_\infty$   
**end for**  
 ▷ Fill Last Layer  $\rightarrow$  super-target  
**for**  $i$  in range( $t$ ) **do**  
    $\mathbf{u}[-i - 2][-1] \leftarrow u_\infty$   
**end for**  
 ▷ Fill  $j$ -th Layer to  $(j + 1)$ -th Layer  
**for**  $j$  in range( $l$ ) **do**  
    $start \leftarrow t * j + 1$   
    $mid \leftarrow t * (j + 1) + 1$   
    $end \leftarrow t * (j + 2) + 1$   
    $\mathbf{u}[start:mid, mid:end] \leftarrow \bar{\mathbf{A}}_{[j,:,:,]}^T$   
**end for**  
 ▷ Get Integral Version of Capacities  
 $\tilde{\mathbf{l}} \leftarrow \text{int}(\gamma * \mathbf{l})$   
 $\tilde{\mathbf{u}} \leftarrow \text{int}(\gamma * \mathbf{u})$   
 ▷ Get Adjacency Matrix  
 $\mathcal{A} \leftarrow \mathbb{I}_{(\mathbf{u} > 0)}$   
 ▷ Get super-source and super-target  
 $ss, st \leftarrow 0, t * (l + 1) + 1$   
**end function**

---

## 3.2 GENERALIZED ATTENTION FLOW

In Generalized Attention Flow, we leverage the attention mechanism for feature attribution by defining a network flow representation of a Transformer or other attention-based model. We assign capacities to the edges of this graph corresponding to information tensor defined in [Sec. 3.1](#). We then solve the maximum flow problem to evaluate the optimal flow passing through any output node (or, more generally, any node in any layer) to any input node. The flow traversing through an input node (token) indicates the importance or attribution of that particular node (token).

To determine the maximum flow from all output nodes to all input nodes, we leverage the concept of multi-commodity flow ([App. A.2](#) and [App. B](#)). This involves the introduction of a super-source node  $ss$  and a super-target node  $st$  with a large capacity  $u_\infty$ . The connectivity between layers and capacities between nodes are established using the information tensors, effectively forming a layered graph ([App. B](#)).

To formalize the generating of the information flow, consider a Transformer with  $l$  attention layers, an input sequence  $X \in \mathbb{R}^{t \times d}$ , and its information tensor  $\tilde{\mathbf{A}} \in \mathbb{R}^{l \times t \times t}$ . We construct the layered attribution graph  $\mathcal{G}$  with its adjacency matrix  $\mathcal{A}$ , its edge-vertex incidence matrix  $\mathcal{B}$ , lower capacity matrix  $\mathbf{l}$  and its integral version  $\tilde{\mathbf{l}}$ , upper capacity matrix  $\mathbf{u}$  and its integral version  $\tilde{\mathbf{u}}$  employing either [Algorithm 1](#) or [Algorithm 2](#). Afterward, we substitute the vectorized version of the obtained matrices into the primal form of [eq. 8](#) to evaluate the desired optimal flow.

To enhance comprehension of [Algorithm 1](#) and [Algorithm 2](#), we explain the process of constructing the layered attribution graph  $\mathcal{G}$  which has an adjacency matrix with shape  $(2 + t * (l + 1), 2 + t * (l + 1))$  and serves as an input for the maximum flow problem. Designating nodes at layer  $\ell \in \{1, \dots, l\}$  and token  $i \in \{1, \dots, t\}$  as  $v_{\ell,i}$ , the guidelines for defining the upper and lower-bound capacities are as follows:

- To connect nodes  $v_{1,i}$  to the super-target node  $v_{st}$ , we define  $\mathbf{u}[0, i] = u_\infty$  for  $1 \leq i \leq t$ .
- The upper-bound capacity from node  $v_{\ell+1,i}$  to node  $v_{\ell,j}$  is defined as  $\mathbf{u}[\mathbf{I}_{i,\ell+1}, \mathbf{I}_{j,\ell}] = \tilde{\mathbf{A}}_{\ell,i,j}$  for  $\ell \in \{1, \dots, l\}$ ,  $i \in \{1, \dots, t\}$ , and  $j \in \{1, \dots, t\}$ , where  $\mathbf{I}_{i,\ell+1} = i + t * \ell$  and  $\mathbf{I}_{j,\ell} = j + t * (\ell - 1)$ .
- To connect the super-source node  $v_{ss}$  to nodes  $v_{l+1,i}$ , we define  $\mathbf{u}[t * l + i, 1 + t * (l + 1)] = u_\infty$  for  $1 \leq i \leq t$ .
- The lower-bound capacity is defined as  $\mathbf{l} = \mathbf{0}$ .

[Fig. 1a](#) and [Fig. 1b](#) illustrate schematic graphs generated using the information tensor  $\tilde{\mathbf{A}} \in \mathbb{R}^{3 \times 3 \times 3}$  with [Algorithm 1](#) and [Algorithm 2](#), respectively. While both algorithms are identically solving the same network flow problem by creating graphs containing a super-source and a super-target, the second algorithm differs from the first in two key aspects. First, in the second graph, the positions of the super-source and super-target are swapped, meaning that the super-source in the first graph becomes the super-target in the second and vice versa. Second, the direction of the edges in the second graph is reversed compared to the first.

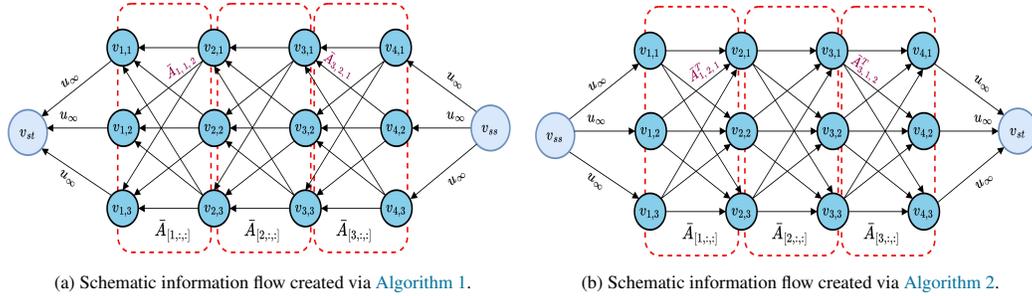


Figure 1: Schematics overview of Generalized Attention Flow created employing [Algorithm 1](#) and [Algorithm 2](#).

### 3.3 NON-UNIQUENESS OF MAXIMUM FLOW

The maximum flow problem lacks strict convexity, meaning it does not necessarily have a unique solution. We found that the maximum flow problem associated with the graphs constructed employing Generalized Attention Flow also fails to yield a unique optimal flow ([App. C](#)).

**Observation 3.1.** *It is straightforward to verify that both [Algorithm 1](#) and [Algorithm 2](#) solve the same maximum flow problem. Therefore, determining the maximum flow in graphs generated by either [Algorithm 1](#) or [Algorithm 2](#) is equivalent and yields the same optimal value. However, it's worth noting that the optimal flows associated with them may not necessarily be equivalent, as explained in [App. C](#).*

282 **Observation 3.2.** *If two distinct feasible solutions, denoted  $f_1$  and  $f_2$ , exist for a linear programming problem,*  
 283 *then any convex combination  $\gamma_1 f_1 + \gamma_2 f_2$  forms another feasible solution. Consequently, the maximum flow*  
 284 *problem can possess an infinite number of feasible solutions. Additionally, due to the non-uniqueness of*  
 285 *optimal flows arising from the maximum flow problem, their projections onto any subset of nodes in the graph*  
 286 *may also not be unique.*

287 **Corollary 3.1.** *Let  $V$  be the set of all nodes in a layered attribution graph  $\mathcal{G}(\mathcal{A}, \mathbf{u}, \mathbf{l}, \mathbf{c}, ss, st)$ , and  $N \subseteq V$ ,*  
 288 *with all nodes in  $N$  chosen from the same layer. Suppose  $\mathbf{f}^*$  is the optimal solution of eq. 8, and for every*  
 289  *$S \subseteq N$ , define the payoff function  $\vartheta(S) := |\mathbf{f}^*(S)| = \sum_{i \in S} |f_{out}(i)|$ , where  $|f_{out}(i)|$  denotes the total outflow*  
 290 *value of node  $i$ . Although [Ethayarajh & Jurafsky \(2021\)](#) claimed that for each node  $i \in N$ ,  $\phi_i(\vartheta) = |\mathbf{f}_{out}^*(i)|$*   
 291 *represents the Shapley value, these feature attributions are non-unique and cannot be considered Shapley*  
 292 *values. In fact, their method for defining feature attributions is not well-defined (Proof in [App. E](#)).*

### 294 3.4 LOG BARRIER REGULARIZATION OF MAXIMUM FLOW

295 To address the non-uniqueness issues in the maximum flow problem, we reformulate the minimum-cost  
 296 circulation problem as follows:

$$297 \arg \min_{\substack{\mathbf{c}^\top \mathbf{f} \\ \mathbf{B}^\top \mathbf{f} = \mathbf{0} \\ \beta(\mathbf{f}) \leq \mathbf{0}}} \mathbf{c}^\top \mathbf{f} \quad (11)$$

298 where  $\beta(\mathbf{f}) = (\mathbf{f} - \mathbf{l})(\mathbf{f} - \mathbf{u})$ . The original problem can, therefore, be approximated using the log barrier  
 299 function as the following optimization problem:

$$300 \arg \min_{\mathbf{B}^\top \mathbf{f} = \mathbf{0}} \mathbf{c}^\top \mathbf{f} + \psi_\mu(\mathbf{f}) \quad (12)$$

301 where the log barrier function is:

$$302 \psi_\mu(\mathbf{f}) = -\mu \sum_{e \in E} \log(-\beta(f_e)) = -\mu \sum_{e \in E} (\log(f_e - l_e) + \log(u_e - f_e)) \quad (13)$$

303 It is evident that, as long as  $\mu > 0$  and our initial solution is feasible, the barrier function guarantees that any  
 304 solution obtained through an iterative minimization scheme, like interior point methods, remains feasible  
 305 ([Bubeck, 2015](#); [Boyd & Vandenberghe, 2004](#); [Mađry, 2019](#)). Furthermore, it can be demonstrated that to  
 306 obtain an  $\varepsilon$ -approximate solution to eq. 11, it suffices to set  $\mu \leq \frac{\varepsilon}{2m}$  and find the optimal solution to the  
 307 corresponding problem in eq. 12 ([Bubeck, 2015](#); [Boyd & Vandenberghe, 2004](#); [Mađry, 2019](#)).

308 Finally, the Hessian of the objective function in eq. 11 at some point  $\mathbf{f}$  is equal to the Hessian  $\nabla^2 \psi_\mu(x)$  of  
 309 the barrier function, which is positive definite (assuming  $\mu > 0$ ). This implies that the objective function is  
 310 strictly convex and, consequently, eq. 12 has a unique feasible solution ([Bubeck, 2015](#); [Boyd & Vandenberghe,](#)  
 311 [2004](#)).

### 312 3.5 AXIOMS OF FEATURE ATTRIBUTIONS

313 In XAI, axioms are core principles that guide the evaluation of explanation methods, ensuring their reliability,  
 314 interpretability, and fairness. These axioms provide standards to measure the effectiveness and compliance of  
 315 explanation techniques. Our proposed methods meet five essential axioms, as demonstrated by the following  
 316 theorem and corollaries.

317 **Definition 3.1 (Shapley values).** *For any value function  $\vartheta : 2^N \mapsto \mathbb{R}$  where  $N = \{1, 2, \dots, n\}$ , Shapley*  
 318 *values  $\phi(\vartheta) \in \mathbb{R}^n$  is computed by averaging the marginal contribution of each feature over all possible*  
 319 *feature combinations as:*

$$\phi_i(\vartheta) = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(n - |S| - 1)!}{n!} (\vartheta(S \cup \{i\}) - \vartheta(S)) \quad (14)$$

Shapley values are the unique explanation that satisfies four fairness-based axioms of efficiency (completeness), symmetry, linearity (additivity), and nullity (Shapley, 1952; Young, 1985) (App. A.3). Initially, a payoff function based on model accuracy was proposed (Lundberg & Lee, 2017); however, since then, various alternative payoff functions have been introduced (Jethani et al., 2022; Sundararajan & Najmi, 2020), each yielding distinct feature importance scores.

**Theorem 3.1 (Log Barrier Regularization of Generalized Attention Flow Outcomes Shapley Values).** *Given a layered attribution graph  $\mathcal{G}(\mathcal{A}, \mathbf{u}, \mathbf{l}, \mathbf{c}, ss, st)$  defined using either Algorithm 1 or Algorithm 2, let  $V$  be the set of all nodes in  $\mathcal{G}$ , and  $N \subseteq V$  such that all nodes in  $N$  are chosen from the same layer. Suppose  $\mathbf{f}^*$  is the optimal unique solution of eq. 12, and for every  $S \subseteq N$ , define the payoff function  $\vartheta(S) := |\mathbf{f}^*(S)| = \sum_{i \in S} |f_{out}(i)|$  where  $|f_{out}(i)|$  is the total outflow value of a node  $i$ . Then, it can be proven that for each node  $i \in N$ ,  $\phi_i(\vartheta) = |f_{out}^*(i)|$  represents the Shapley value (Proof in App. E).*

**Corollary 3.2.** *Theorem 3.1 implies that the feature arbitration obtained by eq. 12 are Shapley values and, consequently, adhere to the axioms of efficiency, symmetry, nullity, linearity.*

## 4 EXPERIMENTS

In this section, we thoroughly evaluate the effectiveness of our methods for sequence classification. While our approach is versatile and applicable to various NLP tasks, including question answering and named entity recognition, which use encoder-only Transformer architectures, this assessment focuses solely on sequence classification.

### 4.1 TRANSFORMER MODELS

Transformer models have demonstrated exceptional performance in various NLP tasks including sequence classification, question answering, and named entity recognition. In our evaluations, we used a specific pre-trained model from the [HuggingFace Hub](#) (Wolf et al., 2020) for each dataset and compared our explanation methods against others to assess their performance (App. F.1).

### 4.2 DATASETS

Our method’s assessment involves sequence classification spanning binary classification tasks on datasets including SST2 (Socher et al., 2013), Amazon Polarity (McAuley & Leskovec, 2013), Yelp Polarity (Zhang et al., 2016), and IMDB (Maas et al., 2011), alongside multi-class classification on the AG News dataset (Zhang et al., 2015). To minimize computational overhead, we conducted experiments on a subset of 5,000 randomly selected samples for the Amazon, Yelp, and IMDB datasets while utilizing the entire test sets for other datasets (App. F.1).

### 4.3 BENCHMARK METHODS

Our experiment compares the methods introduced in Sec. 3.1 with various baseline explanation methods tailored for Transformer models. To evaluate attention-based methods such as RawAtt and Rollout (Abnar & Zuidema, 2020), attention gradient-based methods like Grads, AttGrads (Barkan et al., 2021), CAT, and AttCAT (Qiang et al., 2022), as well as LRP-based methods such as PartialLRP (Voita et al., 2019) and TransAtt (Chefer et al., 2021a), we adapted the repository developed by (Qiang et al., 2022). Additionally, we

376 implemented classical attribution methods such as Integrated Gradient (Sundararajan et al., 2017), KernelShap  
377 (Lundberg & Lee, 2017), and LIME (Ribeiro et al., 2016) leveraging the Captum package (Kokhlikyan et al.,  
378 2020).

#### 380 4.4 EVALUATION METRIC

381  
382 **AOPC:** One of the key evaluation metrics employed is the Area Over the Perturbation Curve (AOPC), a  
383 measure that quantifies the impact of masking top  $k\%$  tokens on the average change in prediction probability  
384 across all test examples. The AOPC is calculated as follows:

$$385 \text{AOPC}(k) = \frac{1}{N} \sum_{i=1}^N p(\hat{y}|\mathbf{x}_i) - p(\hat{y}|\tilde{\mathbf{x}}_i^k) \quad (15)$$

386  
387 where  $N$  is the number of examples,  $\hat{y}$  is the predicted label,  $p(\hat{y}|\cdot)$  is the probability on the predicted  
388 label, and  $\tilde{\mathbf{x}}_i^k$  is constructed by masking the  $k\%$  top-scored tokens from  $\mathbf{x}_i$ . To avoid arbitrary choices for  
389  $k$ , we systematically mask 10%, 20%, ..., 90% of the tokens in order of decreasing saliency, resulting in  
390  $\tilde{\mathbf{x}}_i^{10}, \tilde{\mathbf{x}}_i^{20}, \dots, \tilde{\mathbf{x}}_i^{90}$ .

391  
392 **LOdds:** Log-odds score is calculated by averaging the difference of negative logarithmic probabilities on the  
393 predicted label over all test examples before and after masking  $k\%$  top-scored tokens.

$$394 \text{LOdds}(k) = \frac{1}{N} \sum_{i=1}^N \log \frac{p(\hat{y}|\tilde{\mathbf{x}}_i^k)}{p(\hat{y}|\mathbf{x}_i)} \quad (16)$$

## 398 5 RESULTS

399  
400 We evaluated various explanation methods by masking the top  $k\%$  of tokens across multiple datasets  
401 and measuring their AOPC and LOdds scores, as shown in Tab. 1, which presents average scores for  
402 different  $k$  values. Our findings indicate that the AGF method consistently outperforms others, achieving the  
403 highest AOPC and lowest LOdds scores, effectively identifying and masking the most important tokens for  
404 model predictions. Furthermore, the GF method also exceeds most baseline methods. Evaluation based on  
405 classification metrics also yields consistently similar results (App. D.1 and App. D.2).

406  
407 Additionally, we assessed similar explanation methods by masking the bottom  $k\%$  of tokens across datasets  
408 and measuring AOPC and LOdds scores, detailed in Tab. 2. The AGF method achieved the highest LOdds and  
409 lowest AOPC across most datasets, highlighting its ability to pinpoint important tokens for model predictions,  
410 with the GF method also surpassing many baseline methods in this context.

411 However, the Yelp dataset poses a unique challenge, as our methods do not perform optimally in terms of  
412 AOPC and LOdds metrics. This is likely due to the prevalence of conversational language, slang, and typos in  
413 Yelp reviews, which adversely affect the AGF method’s performance more than others.

## 415 6 LIMITATIONS

416  
417 The primary limitation of our proposed method is the increased running time of the optimization problem  
418 in eq. 12 as the number of tokens grows (Lee & Sidford, 2020; van den Brand et al., 2021). Moreover,  
419 it’s important to note that optimization problems generally cannot be solved in parallel. However, recent  
420 advancements have led to the development of almost-linear time algorithm that solves the optimization  
421 problem described in eq. 12 (Tab. 7). Additionally, we found that the practical runtime of our method is  
422 comparable to other XAI approaches (Tab. 8), indicating that our method is both practical and efficient for  
obtaining feature attributions in AI models using Transformers.

Table 1: AOPC and LOdds scores of all methods in explaining the Transformer-based model across datasets when we mask **top**  $k\%$  tokens. Higher AOPC and lower LOdds are desirable, indicating a strong ability to mark important tokens. Best results are in bold, and differences between AGF and benchmarks are statistically significant according to the ASO test (App. D.3).

Methods	SST2		IMDB		Yelp		Amazon		AG News	
	AOPC $\uparrow$	LOdds $\downarrow$								
RawAtt	0.348	-0.973	0.329	-1.393	0.383	-1.985	0.353	-1.593	0.301	-1.105
Rollout	0.322	-0.887	0.354	-1.456	0.260	-0.987	0.304	-1.326	0.249	-0.983
Grads	0.354	-0.313	0.324	-1.271	0.412	-1.994	0.405	-1.793	0.327	-1.319
AttGrads	0.367	-0.654	0.337	-1.226	0.423	-1.978	0.419	-1.918	0.348	-1.477
CAT	0.369	-1.175	0.332	-1.274	0.417	-1.992	0.381	-1.639	0.325	-1.226
AttCAT	0.405	-1.402	0.371	-1.642	0.431	<b>-2.134</b>	0.427	-2.041	0.387	-1.688
PartialLRP	0.371	-1.171	0.323	-1.321	<b>0.443</b>	-2.018	0.384	-1.945	0.356	-1.627
TransAtt	0.399	-1.286	0.355	-1.513	0.411	-1.473	0.375	-1.875	0.377	-1.318
LIME	0.362	-1.056	0.347	-1.379	0.361	-1.568	0.358	-1.612	0.349	-1.538
KernelShap	0.382	-1.259	0.367	-1.423	0.385	-1.736	0.374	-1.717	0.351	-1.413
IG	0.401	-1.205	0.350	-1.443	0.409	-1.924	0.434	-2.024	0.393	-1.681
AF	0.371	-1.215	0.313	-1.297	0.398	-1.886	0.388	-1.923	0.352	-1.282
GF	0.412	-1.616	0.491	-1.718	0.396	-1.654	0.421	-2.006	0.366	-1.513
AGF	<b>0.427</b>	<b>-1.687</b>	<b>0.498</b>	<b>-1.849</b>	0.429	-1.982	<b>0.439</b>	<b>-2.103</b>	<b>0.398</b>	<b>-1.693</b>

Table 2: AOPC and LOdds scores of all methods in explaining the Transformer-based model across datasets when we mask **bottom**  $k\%$  tokens. Lower AOPC and higher LOdds are desirable, indicating a strong ability to mark important tokens. Best results are in bold, and differences between AGF and benchmarks are statistically significant according to the ASO test (App. D.3).

Methods	SST2		IMDB		Yelp		Amazon		AG News	
	AOPC $\downarrow$	LOdds $\uparrow$								
RawAtt	0.184	-0.693	0.151	-0.471	0.157	-0.747	0.129	-0.281	0.101	-0.427
Rollout	0.221	-0.773	0.123	-0.425	0.169	-0.734	0.171	-0.368	0.117	-0.471
Grads	0.234	-0.776	0.083	-0.203	0.131	-0.641	0.134	-0.254	0.083	-0.390
AttGrads	0.217	-0.713	0.088	-0.243	0.127	-0.603	0.135	-0.266	0.071	-0.351
CAT	0.247	-0.874	0.099	-0.327	0.134	-0.659	0.126	-0.240	0.104	-0.419
AttCAT	0.143	-0.412	0.041	-0.092	<b>0.103</b>	<b>-0.339</b>	0.115	-0.148	0.057	-0.219
PartialLRP	0.163	-0.527	0.057	-0.116	0.116	-0.486	0.167	-0.327	0.056	-0.204
TransAtt	0.148	-0.483	0.045	-0.107	0.123	-0.538	0.113	-0.140	0.049	-0.173
LIME	0.173	-0.603	0.076	-0.141	0.143	-0.687	0.158	-0.263	0.075	-0.372
KernelShap	0.197	-0.729	0.039	-0.084	0.135	-0.645	0.174	-0.351	0.067	-0.219
IG	0.150	-0.532	0.026	-0.064	0.130	-0.617	0.134	-0.241	0.052	-0.191
AF	0.199	-0.747	0.061	-0.148	0.153	-0.689	0.388	-1.923	0.106	-0.402
GF	0.154	-0.497	0.034	-0.079	0.149	-0.654	0.130	-0.267	0.090	-0.313
AGF	<b>0.084</b>	<b>-0.263</b>	<b>0.014</b>	<b>-0.039</b>	0.121	-0.504	<b>0.092</b>	<b>-0.114</b>	<b>0.037</b>	<b>-0.134</b>

## 7 CONCLUSION

In this study, we propose Generalized Attention Flow, an extension of Attention Flow. The core idea behind Generalized Attention Flow is applying the log barrier method to the maximum flow problem, defined by information tensors, to derive feature attributions. By leveraging the log barrier method, we resolve the non-uniqueness issue in optimal flows originating from the maximum flow problem, ensuring that our feature attributions are Shapley values and satisfy efficiency, symmetry, nullity, and linearity axioms. Additionally, we demonstrate that our approach satisfies the axiom of conservation.

Our experiments across multiple datasets indicate that our proposed AGF method generally outperforms other feature attribution methods in most evaluation scenarios. It could be valuable for future research to explore whether alternative definitions of the information tensor could enhance AGF’s effectiveness.

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

## 8 REPRODUCIBILITY

The code used to implement all results presented in this paper is available anonymously in [this repository](#). Comprehensive details for developing the proposed methods can be found in [Sec. 3.1](#), [Sec. 3.2](#), and [App. B](#). Additionally, the pre-trained models, datasets, and further implementation details for our experiments are thoroughly discussed in [Sec. 4](#) and [App. F](#). Proofs supporting our theoretical claims are provided in [App. E](#).

517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563

## REFERENCES

- Samira Abnar and Willem Zuidema. Quantifying Attention Flow in Transformers, May 2020.
- Akshay Agrawal, Robin Verschueren, Steven Diamond, and Stephen Boyd. A rewriting system for convex optimization problems. *Journal of Control and Decision*, 5(1):42–60, 2018.
- Kyriakos Axiotis, Aleksander Madry, and Adrian Vladu. Faster Sparse Minimum Cost Flow by Electrical Flow Localization. In *2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 528–539, Denver, CO, USA, February 2022. IEEE. ISBN 978-1-66542-055-6. doi: 10.1109/FOCS52979.2021.00059.
- Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise Relevance Propagation. *PLOS ONE*, 10(7):e0130140, July 2015. ISSN 1932-6203. doi: 10.1371/journal.pone.0130140.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural Machine Translation by Jointly Learning to Align and Translate, May 2016.
- Oren Barkan, Edan Hauon, Avi Caciularu, Ori Katz, Itzik Malkiel, Omri Armstrong, and Noam Koenigstein. Grad-SAM: Explaining Transformers via Gradient Self-Attention Maps. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pp. 2882–2887, Virtual Event Queensland Australia, October 2021. ACM. ISBN 978-1-4503-8446-9. doi: 10.1145/3459637.3482126.
- Stephen P. Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, Cambridge, UK ; New York, 2004. ISBN 978-0-521-83378-3.
- Sébastien Bubeck. *Convex Optimization: Algorithms and Complexity*, November 2015.
- Hila Chefer, Shir Gur, and Lior Wolf. Transformer Interpretability Beyond Attention Visualization. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 782–791, Nashville, TN, USA, June 2021a. IEEE. ISBN 978-1-66544-509-2. doi: 10.1109/CVPR46437.2021.00084.
- Hila Chefer, Shir Gur, and Lior Wolf. Generic Attention-model Explainability for Interpreting Bi-Modal and Encoder-Decoder Transformers, March 2021b.
- Hanjie Chen, Guangtao Zheng, and Yangfeng Ji. Generating Hierarchical Explanations on Text Classification via Feature Interaction Detection, May 2020.
- Li Chen, Rasmus Kyng, Yang P. Liu, Simon Meierhans, and Maximilian Probst Gutenberg. Almost-Linear Time Algorithms for Incremental Graphs: Cycle Detection, SCCs, Shortest Path, and Minimum-Cost Flow, November 2023a.
- Lu Chen, Siyu Lou, Keyan Zhang, Jin Huang, and Quanshi Zhang. HarsanyiNet: Computing Accurate Shapley Values in a Single Forward Propagation, December 2023b.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. What Does BERT Look at? An Analysis of BERT’s Attention. In Tal Linzen, Grzegorz Chrupała, Yonatan Belinkov, and Dieuwke Hupkes (eds.), *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pp. 276–286, Florence, Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-4828.
- Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms, Third Edition*. The MIT Press, 3rd edition, August 2009. ISBN 978-0-262-03384-8.
- E. del Barrio, J. A. Cuesta-Albertos, and C. Matrán. An optimal transportation approach for assessing almost stochastic order, May 2017.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, May 2019.

564 Steven Diamond and Stephen Boyd. CVXPY: A Python-embedded modeling language for convex optimization. *Journal*  
565 *of Machine Learning Research*, 17(83):1–5, 2016.

566

567 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa  
568 Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An Image is Worth  
569 16x16 Words: Transformers for Image Recognition at Scale. In *International Conference on Learning Representations*,  
570 October 2020.

571 Rotem Dror, Segev Shlomov, and Roi Reichart. Deep Dominance - How to Properly Compare Deep Neural Models. In  
572 Anna Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Annual Meeting of the Association*  
573 *for Computational Linguistics*, pp. 2773–2785, Florence, Italy, July 2019. Association for Computational Linguistics.  
574 doi: 10.18653/v1/P19-1266.

575 Kawin Ethayarajh and Dan Jurafsky. Attention Flows are Shapley Value Explanations, May 2021.

576 Javier Ferrando and Elena Voita. Information Flow Routes: Automatically Interpreting Language Models at Scale, April  
577 2024.

578 Sarthak Jain and Byron C. Wallace. Attention is not Explanation. In Jill Burstein, Christy Doran, and Tamar Solorio  
579 (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational*  
580 *Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 3543–3556, Minneapolis,  
581 Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1357.

582 Neil Jethani, Mukund Sudarshan, Ian Covert, Su-In Lee, and Rajesh Ranganath. FastSHAP: Real-Time Shapley Value  
583 Estimation, March 2022.

584

585 Goro Kobayashi, Tatsuki Kuribayashi, Sho Yokoi, and Kentaro Inui. Incorporating Residual and Normalization Layers  
586 into Analysis of Masked Language Models. In *Proceedings of the 2021 Conference on Empirical Methods in*  
587 *Natural Language Processing*, pp. 4547–4568, Online and Punta Cana, Dominican Republic, 2021. Association for  
588 Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.373.

589 Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds, Alexander  
590 Melnikov, Natalia Kliushkina, Carlos Araya, Siqi Yan, and Orion Reblitz-Richardson. Captum: A unified and generic  
591 model interpretability library for PyTorch, September 2020.

592 Yin Tat Lee and Aaron Sidford. Path Finding Methods for Linear Programming: Solving Linear Programs in  
593  $\epsilon$  Iterations and Faster Algorithms for Maximum Flow. In *2014 IEEE 55th Annual Symposium*  
594 *on Foundations of Computer Science*, pp. 424–433, Philadelphia, PA, USA, October 2014. IEEE. ISBN 978-1-4799-  
595 6517-5. doi: 10.1109/FOCS.2014.52.

596 Yin Tat Lee and Aaron Sidford. Solving Linear Programs with  $\sqrt{\text{rank}}$  Linear System Solves, August 2020.

597

598 Scott Lundberg and Su-In Lee. A Unified Approach to Interpreting Model Predictions, November 2017.

599 Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning Word  
600 Vectors for Sentiment Analysis. In Dekang Lin, Yuji Matsumoto, and Rada Mihalcea (eds.), *Proceedings of the 49th*  
601 *Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pp. 142–150,  
602 Portland, Oregon, USA, June 2011. Association for Computational Linguistics.

603 Aleksander Mądry. GRADIENTS AND FLOWS: CONTINUOUS OPTIMIZATION APPROACHES TO THE MAX-  
604 IMUM FLOW PROBLEM. In *Proceedings of the International Congress of Mathematicians (ICM 2018)*, pp.  
605 3361–3387, Rio de Janeiro, Brazil, May 2019. WORLD SCIENTIFIC. ISBN 978-981-327-287-3 978-981-327-288-0.  
606 doi: 10.1142/9789813272880\_0185.

607 Julian McAuley and Jure Leskovec. Hidden factors and hidden topics: Understanding rating dimensions with review text.  
608 In *Proceedings of the 7th ACM Conference on Recommender Systems*, pp. 165–172, Hong Kong China, October 2013.  
609 ACM. ISBN 978-1-4503-2409-0. doi: 10.1145/2507157.2507163.

610

611 Ali Modarressi, Mohsen Fayyaz, Yadollah Yaghoobzadeh, and Mohammad Taher Pilehvar. GlobEnc: Quantifying Global  
612 Token Attribution by Incorporating the Whole Encoder Layer in Transformers. In Marine Carpuat, Marie-Catherine de  
613 Marneffe, and Ivan Vladimir Meza Ruiz (eds.), *Proceedings of the 2022 Conference of the North American Chapter of  
614 the Association for Computational Linguistics: Human Language Technologies*, pp. 258–271, Seattle, United States,  
615 July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.19.

616 Ali Modarressi, Mohsen Fayyaz, Ehsan Aghazadeh, Yadollah Yaghoobzadeh, and Mohammad Taher Pilehvar. DecompX:  
617 Explaining Transformers Decisions by Propagating Token Decomposition. In Anna Rogers, Jordan Boyd-Graber,  
618 and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics  
619 (Volume 1: Long Papers)*, pp. 2649–2664, Toronto, Canada, July 2023. Association for Computational Linguistics. doi:  
620 10.18653/v1/2023.acl-long.149.

621 Aaron Mueller. Missed Causes and Ambiguous Effects: Counterfactuals Pose Challenges for Interpreting Neural  
622 Networks. 2024. doi: 10.48550/ARXIV.2407.04690.

623 Dong Nguyen. Comparing Automatic and Human Evaluation of Local Explanations for Text Classification. In Marilyn  
624 Walker, Heng Ji, and Amanda Stent (eds.), *Proceedings of the 2018 Conference of the North American Chapter of the  
625 Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 1069–1078,  
626 New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1097.

627 Yao Qiang, Deng Pan, Chengyin Li, Xin Li, Rhongho Jang, and Dongxiao Zhu. AttCAT: Explaining Transformers via  
628 Attentive Class Activation Tokens. In *Advances in Neural Information Processing Systems*, October 2022.

629 Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "Why Should I Trust You?": Explaining the Predictions of Any  
630 Classifier, August 2016.

631 Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. DistilBERT, a distilled version of BERT: Smaller,  
632 faster, cheaper and lighter, February 2020.

633 Sofia Serrano and Noah A. Smith. Is Attention Interpretable? In *Proceedings of the 57th Annual Meeting of the  
634 Association for Computational Linguistics*, pp. 2931–2951, Florence, Italy, 2019. Association for Computational  
635 Linguistics. doi: 10.18653/v1/P19-1282.

636 Lloyd S. Shapley. A Value for N-Person Games. Technical report, RAND Corporation, March 1952.

637 Avanti Shrikumar, Jocelin Su, and Anshul Kundaje. Computationally Efficient Measures of Internal Neuron Importance,  
638 July 2018.

639 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts.  
640 Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank. In David Yarowsky, Timothy  
641 Baldwin, Anna Korhonen, Karen Livescu, and Steven Bethard (eds.), *Proceedings of the 2013 Conference on Empirical  
642 Methods in Natural Language Processing*, pp. 1631–1642, Seattle, Washington, USA, October 2013. Association for  
643 Computational Linguistics.

644 Mukund Sundararajan and Amir Najmi. The Many Shapley Values for Model Explanation. In *Proceedings of the 37th  
645 International Conference on Machine Learning*, pp. 9269–9278. PMLR, November 2020.

646 Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic Attribution for Deep Networks, June 2017.

647 Dennis Ulmer, Christian Hardmeier, and Jes Frellsen. Deep-significance - Easy and Meaningful Statistical Significance  
648 Testing in the Age of Neural Networks, April 2022.

649 Jan van den Brand, Yu Gao, Arun Jambulapati, Yin Tat Lee, Yang P. Liu, Richard Peng, and Aaron Sidford. Faster  
650 Maxflow via Improved Dynamic Spectral Vertex Sparsifiers, December 2021.

651 Jan Van Den Brand, Yin Tat Lee, Yang P. Liu, Thatchaphol Saranurak, Aaron Sidford, Zhao Song, and Di Wang. Minimum  
652 cost flows, MDPs, and  $\ell_1$ -regression in nearly linear time for dense instances. In *Proceedings of the 53rd Annual ACM  
653 SIGACT Symposium on Theory of Computing*, pp. 859–869, Virtual Italy, June 2021. ACM. ISBN 978-1-4503-8053-9.  
654  
655  
656  
657

658           doi: 10.1145/3406325.3451108.  
659  
660 Jan van den Brand, Yang P. Liu, and Aaron Sidford. Dynamic Maxflow via Dynamic Interior Point Methods. In  
661 *Proceedings of the 55th Annual ACM Symposium on Theory of Computing, STOC 2023*, pp. 1215–1228, New York, NY,  
662 USA, June 2023. Association for Computing Machinery. ISBN 978-1-4503-9913-5. doi: 10.1145/3564246.3585135.

663 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia  
664 Polosukhin. Attention is All you Need. In *Advances in Neural Information Processing Systems*, volume 30. Curran  
665 Associates, Inc., 2017.

666 Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing Multi-Head Self-Attention:  
667 Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned. In Anna Korhonen, David Traum, and Lluís  
668 Màrquez (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp.  
669 5797–5808, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1580.

670 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim  
671 Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite,  
672 Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush.  
673 HuggingFace’s Transformers: State-of-the-art Natural Language Processing, July 2020.

674 H. P. Young. Monotonic solutions of cooperative games. *International Journal of Game Theory*, 14(2):65–72, June 1985.  
675 ISSN 1432-1270. doi: 10.1007/BF01769885.

676 Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level Convolutional Networks for Text Classification. In *Advances*  
677 *in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015.

678  
679 Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level Convolutional Networks for Text Classification, April 2016.  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704

705 A PRELIMINARIES

706  
707 A.1 MAXIMUM FLOW

708  
709 **Definition A.1 (Network Flow).** Given a network  $G = (V, E, s, t, \mathbf{u})$ , where  $s$  and  $t$  are the source and  
710 target nodes respectively and  $u_{ij}$  is the capacity for the edge  $(i, j) \in E$ , a flow is characterized as a function  
711  $f : E \rightarrow \mathbb{R}^{\geq 0}$  s.t.

$$712 \quad f_{ij} \leq u_{ij} \quad \forall (i, j) \in E \quad (\text{capacity constraints})$$

$$713 \quad \sum_{j:(i,j) \in E} f_{ij} - \sum_{j:(j,i) \in E} f_{ji} = 0, \quad \forall i \in V, i \neq s, t \quad (\text{flow conservation constraints}) \quad (17)$$

714  
715  
716 We define  $|f_{\text{out}}(i)|$  to be the total outflow value of a node  $i$  and  $|f_{\text{in}}(i)|$  to be the total inflow value of a node  $i$ .  
717 For a given set  $K \subseteq V$  of nodes, we define  $|f(K)| = \sum_{i \in K} |f_{\text{out}}(i)|$  for every flow  $f$ . The value of a flow in  
718 a given network  $G = (V, E, s, t, \mathbf{u})$  is denoted as  $|f| = \sum_{v:(s,v)} f_{sv} - \sum_{v:(v,s)} f_{vs} = |f_{\text{out}}(s)| - |f_{\text{in}}(s)|$ ,  
719 and a maximum flow is identified as a feasible flow with the highest attainable value.  
720

721  
722 A.2 MULTI-COMMODITY MAXIMUM FLOW

723 The multi-commodity flow problem is an important variant of the maximum flow problem. This problem  
724 involves multiple source-sink pairs, unlike the standard maximum flow problem, which only considers one  
725 source and one sink. The goal is to find multiple optimal flows, denoted by  $f^1(\cdot, \cdot), \dots, f^r(\cdot, \cdot)$ , where each  
726  $f^k(\cdot, \cdot)$  represents a feasible flow from the source  $s_k$  to the sink  $t_k$ . The objective is to ensure that all capacity  
727 constraints are satisfied, which are represented by the equation:

$$728 \quad \sum_{k=1}^r f^k(i, j) \leq u(i, j) \quad \forall (i, j) \in E \quad (18)$$

729 Such a flow is known as a "multi-commodity" flow. A multi-commodity maximum flow problem is to  
730 maximize the function  $\sum_{k=1}^r \sum_{v:(v,s_k)} f^k(s_k, v)$ .

731 To solve the problem of multi-commodity maximum flow, we can simplify it by transforming it into a standard  
732 maximum flow problem. This can be achieved by introducing two new nodes, a "super-source" node  $ss$  and a  
733 "super-target" node  $st$ . The "super-source" node  $ss$  should be connected to all the original sources  $s_i$  through  
734 edges of finite capacities, while the "super-target" node  $st$  should be connected to all the original sinks  $t_i$  with  
735 edges of finite capacities:

- 736 • Each outgoing edge from the "super-source" node  $ss$  to each source node  $s_i$  gets assigned a capacity  
737 that is equal to the total capacity of the outgoing edges from the source node  $s_i$ .
- 738 • Each incoming edge from an original "super-target" node  $st$  to each sink node  $t_i$  gets assigned a  
739 capacity that is equal to the total capacity of the incoming edges to the sink node  $t_i$ .

740 It is easy to demonstrate that the maximum flow from  $ss$  to  $st$  is equivalent to the maximum sum of flows in a  
741 feasible multi-commodity flow in the original network.  
742

743  
744 A.3 SHAPLEY VALUES

745 The Shapley value, introduced by Shapley (1952), concerns the cooperative game in the coalitional form  
746  $(N, \vartheta)$ , where  $N$  is a set of  $n$  players and  $\vartheta : 2^N \rightarrow \mathbb{R}$  with  $\vartheta(\emptyset) = 0$  is the characteristic (payoff)  
747 function. In the game, the marginal contribution of the player  $i$  to any coalition  $S$  with  $i \notin S$  is considered as  
748

$\vartheta(S \cup i) - \vartheta(S)$ . These Shapley values are the only constructs that jointly satisfy the efficiency, symmetry, nullity, and additivity axioms (Shapley, 1952; Young, 1985):

**Efficiency:** The Shapley values must add up to the total value of the game, which means  $\sum_{i \in N} \phi_i(\vartheta) = \vartheta(N)$ .

**Symmetry:** If two players are equal in their contributions to any coalition, they should receive the same Shapley value. Mathematically, if  $\vartheta(S \cup \{i\}) = \vartheta(S \cup \{j\})$  for all  $S \subseteq N \setminus \{i, j\}$ , then  $\phi_i(\vartheta) = \phi_j(\vartheta)$ .

**Nullity (Dummy):** If a player has no impact on any coalition, their Shapley value should be zero. Mathematically, if  $\vartheta(S \cup \{i\}) = \vartheta(S)$  for all  $S \subseteq N \setminus \{i\}$ , then  $\phi_i(\vartheta) = 0$ .

**Linearity:** If the game  $\vartheta(\cdot)$  is a linear combination of two games  $\vartheta_1(\cdot), \vartheta_2(\cdot)$  for all  $S \subseteq N$ , i.e.  $\vartheta(S) = \vartheta_1(S) + \vartheta_2(S)$  and  $(c \cdot \vartheta)(S) = c \cdot \vartheta(S), \forall c \in \mathbb{R}$ , then the Shapley value in the game  $\vartheta$  is also a linear combination of that in the games  $\vartheta_1$  and  $\vartheta_2$ , i.e.  $\forall i \in N, \phi_i(\vartheta) = \phi_i(\vartheta_1) + \phi_i(\vartheta_2)$  and  $\phi_i(c \cdot \vartheta) = c \cdot \phi_i(\vartheta)$ .

Considering these axioms, the attribution of a player  $j$  is uniquely given by (Shapley, 1952; Young, 1985):

$$\phi_j(\vartheta) = \sum_{S \subseteq N \setminus \{j\}} \frac{|S|!(n - |S| - 1)!}{n!} (\vartheta(S \cup \{j\}) - \vartheta(S)) \quad (19)$$

where the difference  $\vartheta(S \cup \{j\}) - \vartheta(S)$  represents the  $i$ -th feature's contribution to the subset  $S$ , and the summation represents a weighted average across all subsets that do not include  $i$ .

Initially, a payoff function based on model accuracy was suggested (Lundberg & Lee, 2017). Since then, various alternative coalition functions have been proposed (Jethani et al., 2022; Sundararajan & Najmi, 2020), each resulting in a different feature importance score. Many of these alternative approaches are widely used and have been shown to outperform the basic SHAP method (Lundberg & Lee, 2017) in empirical studies.

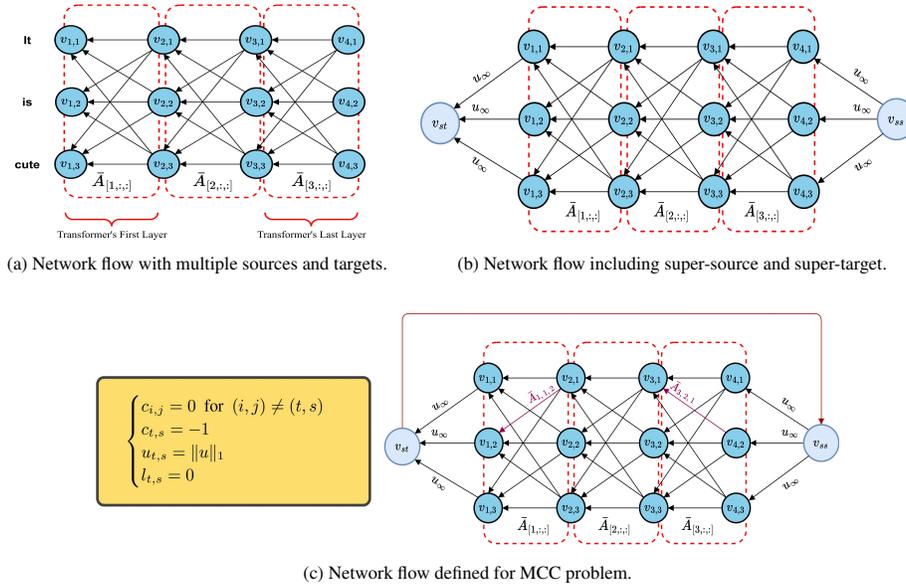


Figure 2: Initial network flow to be used in our proposed method, the multi-commodity flow with multiple sources and targets, and the network flow for MCC problems.

## B NETWORK FLOW GENERATION

Fig. 2 will detail the process of defining a graph network and its parameters using Algorithm 1 for use in our proposed method. It is worth noting that the same method can be used for Algorithm 2. To solve the maximum flow or MCC problem within this graph network, we must compute network flow with multiple sources and targets, assigning all nodes in the first and last layers of transformers as sources and targets, respectively (Fig. 2a).

To solve this problem, we leverage the concept of multi-commodity flow (multiple-sources multiple-targets maximum flow) by introducing a super-source node  $ss$  and a super-target node  $st$  (Fig. 2b). To define the upper-bound and lower-bound capacities of this new graph network, we utilize the procedure defined in Sec. 3.2. In the last step, we add a new edge from the super-target node  $st$  to the super-source node  $ss$  and define the cost vector, upper-bound capacities, and lower-band capacities according to Fig. 2c. Subsequently, we can input all derived parameters into eq. 12, solve the optimization problem, and evaluate feature attributions.

## C NON-UNIQUENESS OF MAXIMUM FLOW

Fig. 3 visually describes our proposed approach for computing feature attributions. Using maximum flow to derive these attributions produces a convex set containing all optimal flows, which makes it unsuitable as a feature attribution technique. In contrast, our proposed approach, which utilizes the log barrier method, generate a unique optimal flow and provides an interpretable set of feature attributions.

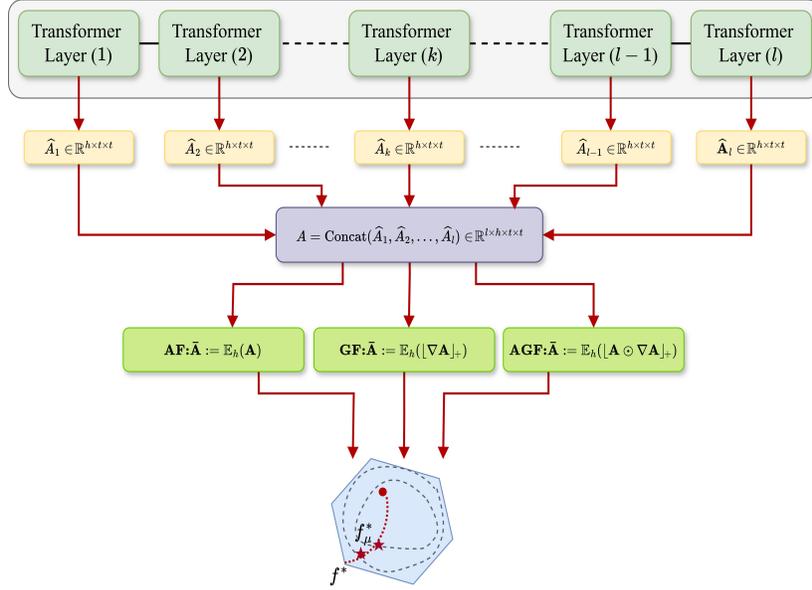
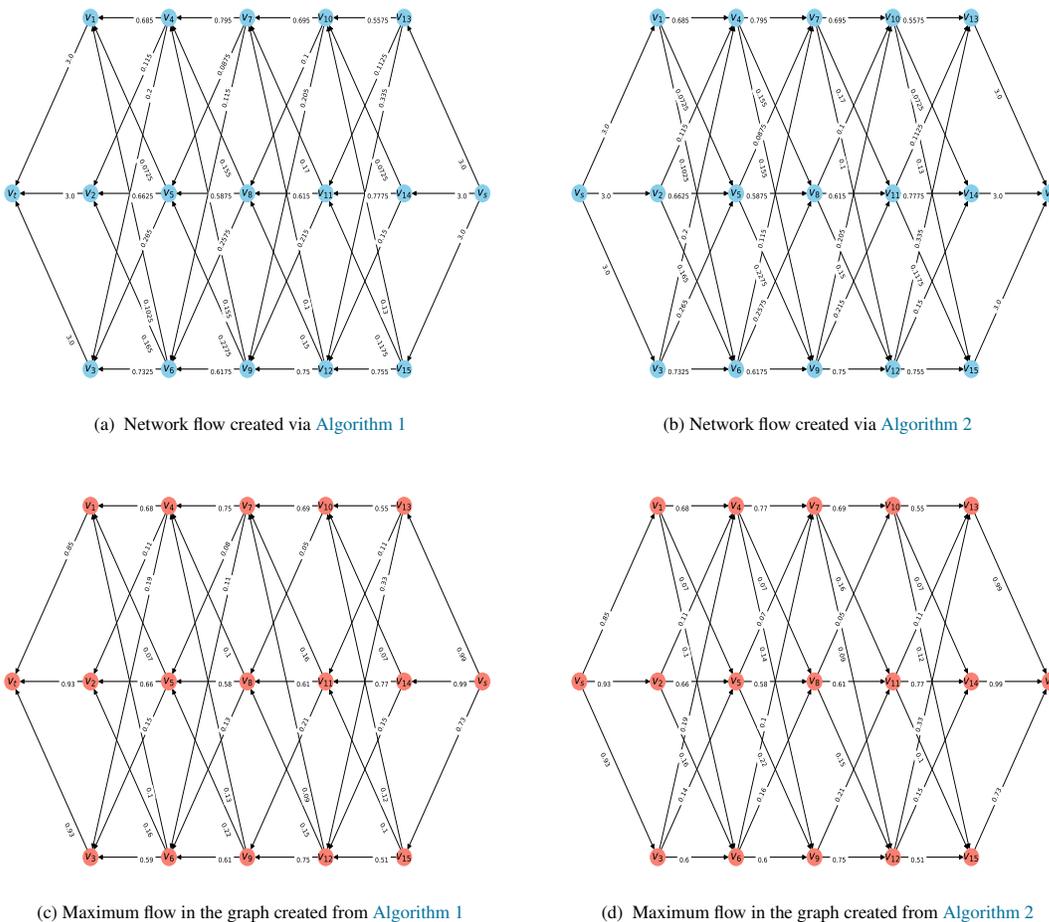


Figure 3: Overview of how the proposed method evaluates the unique optimal flow computed using the log barrier method, attention weights, and their gradients in Transformers.

Fig. 4 shows the capacities and optimal flows obtained by solving maximum flows on the network, constructed with the synthetic information tensor  $\hat{\mathbf{A}} \in \mathbb{R}^{4 \times 3 \times 3}$  as input, using Algorithm 1 and Algorithm 2. While the maximum flows are the same for both algorithms, their optimal flows differ. Notably, significant differences in

846 flows between node pairs  $\{v_4, v_8\}$  and  $\{v_3, v_5\}$  are visible in Fig. 4c and Fig. 4d. Additionally, we evaluated  
 847 the maximum flow and its optimal flow generated by both algorithms across various combinations of token  
 848 numbers  $t$  and Transformer layers  $l$ . Our findings indicate that the optimal flows from the two algorithms do  
 849 not coincide in any scenario.  
 850  
 851

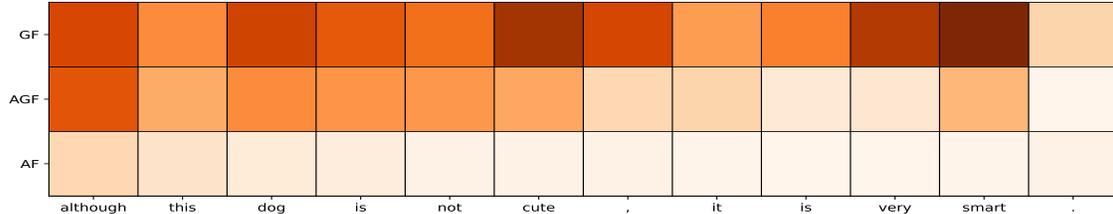


852  
 853  
 854  
 855  
 856  
 857  
 858  
 859  
 860  
 861  
 862  
 863  
 864  
 865  
 866  
 867  
 868  
 869  
 870  
 871  
 872  
 873  
 874  
 875  
 876  
 877  
 878  
 879  
 880  
 881  
 882 Figure 4: Network flow, maximum flow, and residual flow created by Algorithm 1 and Algorithm 2. The optimal flows  
 883 and residual flows evaluated using Algorithm 1 and Algorithm 2 are different.  
 884

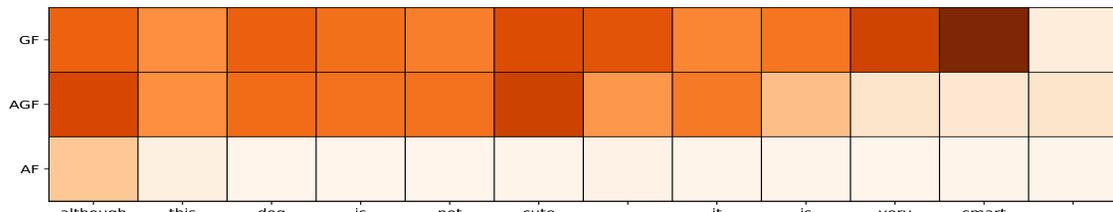
885 Fig. 5a and Fig. 5b display the normalized feature attributions evaluated across three information tensors  
 886 introduced in Sec. 3.1 for sentiment analysis of the sentence "although this dog is not cute, it is very smart."  
 887 employing both Algorithm 1 and Algorithm 2. Across each of the three information tensors, the resulting  
 888 optimal flows and their corresponding normalized attributions differ depending on whether Algorithm 1 or  
 889 Algorithm 2 is used.

890 The layer-wise normalized feature attributions, obtained through the same process, are displayed in Fig. 6.  
 891 For each information tensor type and layer, the resulting optimal flows and their normalized attributions differ  
 892 based on whether Algorithm 1 or Algorithm 2 is utilized. We also computed the optimal flow and its feature

893 attributions for various input sentences using both algorithms for each of the information tensors AF, GF,  
 894 and AGF. Our findings reveal that the optimal flows and corresponding feature attributions generated by  
 895 [Algorithm 1](#) and [Algorithm 2](#) differ for all sentences.  
 896

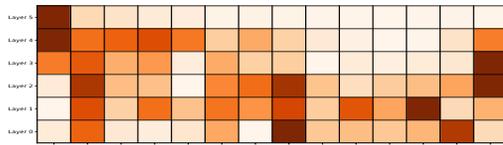


(a) Normalized feature attributions for Transformer's input layer evaluated by [Algorithm 1](#) for different information tensors.

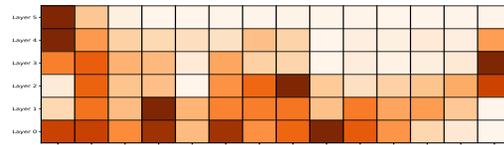


(b) Normalized feature attributions for Transformer's input layer evaluated by [Algorithm 2](#) different information tensors.

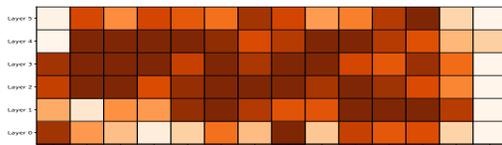
Figure 5: Normalized feature attributions for Transformer's input layer and different information tensors.



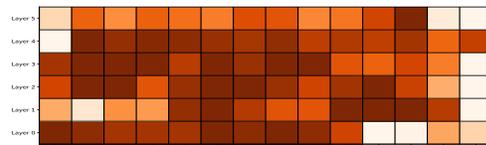
(a) AF method:[Algorithm 1](#)



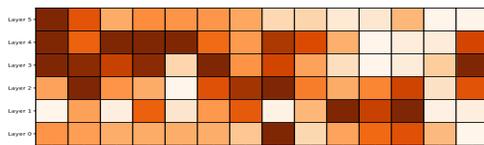
(b) AF method:[Algorithm 2](#).



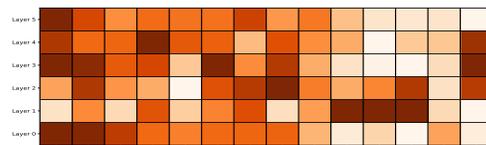
(c) GF method:[Algorithm 1](#).



(d) GF method:[Algorithm 2](#).



(e) AGF method:[Algorithm 1](#).



(f) AGF method:[Algorithm 2](#).

Figure 6: Normalized feature attributions for all Transformer layers evaluated by [Algorithm 1](#) and [Algorithm 2](#).

## D RESULTS

### D.1 QUALITATIVE VISUALIZATIONS

This section visually examines the feature attributions derived from our proposed methods, applied to information tensors as defined in Sec. 3.1. Fig. 7 illustrates feature attributions obtained from our proposed methods applied to two graphs generated by either Algorithm 1 or Algorithm 2. Remarkably, our approaches consistently yield identical results for both graphs. The outcomes vividly demonstrate the superiority of AGF over AF and GF, offering more insightful and reasonable feature attributions. Specifically, both AGF and GF effectively highlight the importance of tokens like 'smart' and 'cute', while assigning lower values to less significant tokens such as 'this', 'it', and 'and'. In contrast, AF fails to capture the expected feature attribution of 'smart' and tends to produce an almost uniform distribution for feature attributions.

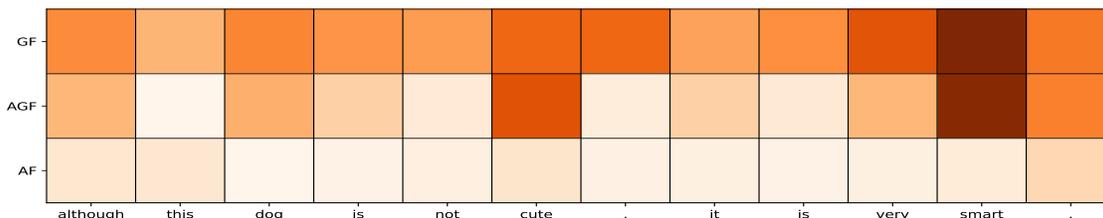


Figure 7: Visualizations of the feature attributions generated by running our proposed method on the three introduced information tensors on the showcase example.

### D.2 ADDITIONAL RESULTS

Fig. 8 presents a detailed comparison of the performance dynamics of various feature attribution methods under different corruption rates across three distinct datasets: IMDB, Amazon, and Yelp. The efficacy of these methods is evaluated using two key metrics, AOPC and LOdds, which measure how well each method identifies important tokens that influence model predictions. Notably, our proposed AGF method consistently outperforms the other techniques, maintaining the highest average AOPC and LOdds scores across a range of corruption levels, particularly for both the IMDB and Amazon datasets. This consistent superiority highlights the robustness of AGF in pinpointing the most important tokens, which significantly affect the model's decision-making process.

We also evaluated various explanation methods by analyzing their performance on classification metrics, with results summarized in Tab. 3, which details the average Accuracy, F1, Precision, and Recall scores across multiple  $k$  values. On the SST2 dataset, our proposed methods AGF and GD, in conjunction with KernelShap, achieved the highest overall performance. For the IMDB dataset, AGF and GF, along with Integrated Gradients (IG), showed significant improvement over other methods. In the Amazon dataset, AGF and GF, when combined with TransAtt, outperformed all competitors. For the AG News dataset, AGF and AF, paired with AttGrads, demonstrated superior performance. The consistently strong performance of AGF across diverse datasets and evaluation metrics highlight its versatility and effectiveness in accurately identifying important tokens, reinforcing its reliability as a feature attribution method in NLP models.

However, the Yelp dataset poses a distinct challenge where our proposed methods, including AGF, do not consistently achieve optimal results across all evaluation metrics. This performance discrepancy is likely due to the unique characteristics inherent to the Yelp dataset, which frequently contains a higher concentration of informal language, colloquialisms, and typographical errors. The prevalence of such linguistic noise in Yelp reviews is notably higher compared to other datasets like IMDB or Amazon. These textual irregularities

introduce complexity that AGF, with its current configuration, may be less equipped to handle effectively. Consequently, AGF’s performance suffers, as it appears to have a lower tolerance for handling noisy or non-standard text inputs, which compromise its ability to accurately attribute features and identify the most influential tokens within these reviews.

Fig. 9 compare the feature attribution methods evaluated on the aforementioned sentence using a model trained on SST2. In all instances, the feature attribution methods predict positive sentiment for the showcased example. Our methods, AGF and GF, effectively capture the most important tokens, such as ‘cute’ and ‘smart’ (indicated in dark orange shading), which significantly contribute to the positive sentiment prediction. Some other methods, including Grads, LIME, RawAtt, and PartialLRP, also exhibit some capability in identifying important tokens. However, certain methods like AF, AttCAT, CAT, Rollout, KernelShap, and IG struggle to correctly identify important tokens.

Table 3: The average of F1, Accuracy, Precision, and Recall scores of all methods in explaining the Transformer-based model on each dataset when we mask **top**  $k\%$  tokens. Lower scores are desirable for all metrics (indicated by  $\downarrow$ ), indicating a strong ability to mark important tokens.

Methods	SST2				IMDB				Yelp				Amazon				AG News				
	F1 $\downarrow$	Acc $\downarrow$	Prec $\downarrow$	Rec $\downarrow$	F1 $\downarrow$	Acc $\downarrow$	Prec $\downarrow$	Rec $\downarrow$	F1 $\downarrow$	Acc $\downarrow$	Prec $\downarrow$	Rec $\downarrow$	F1 $\downarrow$	Acc $\downarrow$	Prec $\downarrow$	Rec $\downarrow$	F1 $\downarrow$	Acc $\downarrow$	Prec $\downarrow$	Rec $\downarrow$	
RawAtt	0.75	0.75	0.72	0.79	0.69	0.67	0.70	0.69	0.68	0.72	0.74	0.63	0.67	0.68	0.67	0.66	0.65	0.68	0.68	0.68	0.68
Rollout	0.81	0.82	0.80	0.82	0.74	0.67	0.65	0.84	0.80	0.83	0.88	0.74	0.71	0.73	0.71	0.72	0.61	0.63	0.64	0.63	0.63
Grads	0.78	0.75	0.72	0.79	0.69	0.67	0.70	0.69	0.68	0.72	0.74	0.63	0.67	0.68	0.67	0.66	0.65	0.68	0.68	0.68	0.68
AttGrads	0.78	0.78	0.75	0.82	0.76	0.75	0.78	0.76	0.91	0.91	0.89	0.93	0.79	0.82	0.80	0.78	0.58	0.61	0.60	0.60	0.60
CAT	0.68	0.65	0.61	0.76	0.56	0.49	0.51	0.65	0.70	0.70	0.68	0.73	0.68	0.64	0.67	0.70	0.63	0.64	0.64	0.64	0.64
AttCAT	0.68	0.65	0.62	0.76	0.57	0.48	0.50	0.64	0.67	0.66	0.65	0.70	0.67	0.62	0.66	0.68	0.64	0.64	0.65	0.64	0.64
PartialLRP	0.75	0.75	0.71	0.78	0.66	0.65	0.67	0.67	0.65	0.70	0.71	0.61	0.65	0.66	0.65	0.64	0.65	0.68	0.68	0.68	0.68
TransAtt	0.73	0.72	0.69	0.76	0.61	0.58	0.60	0.62	0.62	0.66	0.66	0.60	0.62	0.63	0.63	0.61	0.63	0.66	0.66	0.66	0.66
LIME	0.61	0.63	0.62	0.63	0.55	0.55	0.55	0.55	0.72	0.73	0.72	0.73	0.72	0.73	0.72	0.73	0.67	0.67	0.68	0.67	0.67
KernelShap	0.53	0.52	0.53	0.53	0.67	0.69	0.68	0.69	0.77	0.78	0.77	0.78	0.67	0.68	0.67	0.68	0.74	0.74	0.75	0.74	0.74
IG	0.56	0.57	0.56	0.57	0.48	0.50	0.48	0.50	0.72	0.72	0.72	0.72	0.73	0.74	0.73	0.74	0.67	0.68	0.67	0.67	0.67
AF	0.72	0.72	0.71	0.71	0.65	0.68	0.70	0.68	0.71	0.71	0.71	0.70	0.69	0.71	0.71	0.71	0.64	0.62	0.65	0.64	0.64
GF	0.56	0.57	0.56	0.56	0.45	0.48	0.45	0.48	0.70	0.71	0.70	0.71	0.65	0.67	0.66	0.67	0.67	0.68	0.67	0.67	0.67
AGF	0.54	0.52	0.54	0.54	0.46	0.47	0.47	0.47	0.70	0.72	0.70	0.70	0.67	0.67	0.67	0.64	0.63	0.65	0.64	0.64	0.64

### D.3 STATISTICAL SIGNIFICANCE TEST

To implement a statistical significance test, we employ the ASO (Almost Stochastic Order) method (Ulmer et al., 2022; Dror et al., 2019; del Barrio et al., 2017), which compares the cumulative distribution functions (CDFs) of two score distributions to determine stochastic dominance. Notably, ASO imposes no assumptions about score distributions, making it applicable to any metric where higher scores indicate better performance.

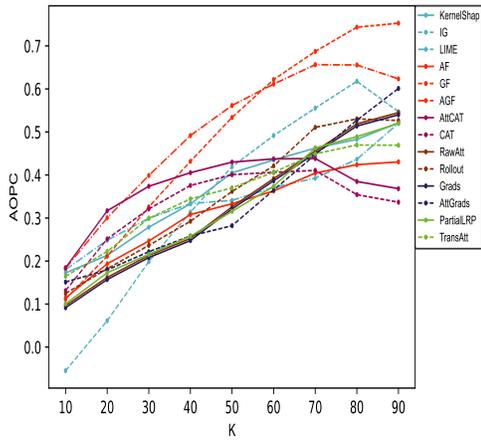
When comparing model  $A$  with model  $B$  using the ASO method, we obtain the value  $\epsilon_{\min}$ , which is an upper bound on the violation of stochastic order. If  $\epsilon_{\min} \leq \tau$  (with  $\tau \leq 0.5$ ), model  $A$  is considered stochastically dominant over model  $B$ , implying superiority. This value can also be interpreted as a confidence score; a lower  $\epsilon_{\min}$  suggests greater confidence in model  $A$ ’s superiority. The null hypothesis for ASO is defined as:

$$H_0 : \epsilon_{\min} \geq \tau \quad (20)$$

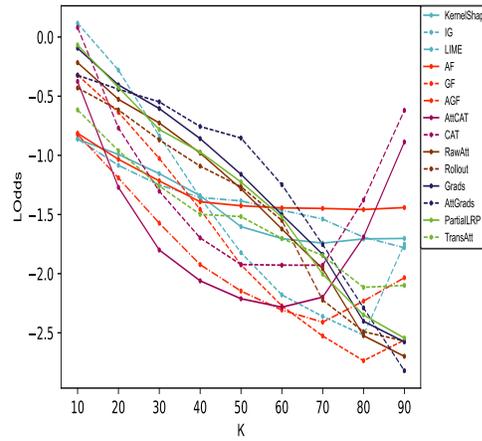
where the significance level  $\alpha$  is an input parameter that influences  $\epsilon_{\min}$ .

In this study, we conduct 500 independent runs per method to perform comprehensive statistical tests, comparing the AOPC and LOdds metrics between our best-performing proposed method, AGF, and the top benchmark methods outlined in Tab. 1 and Tab. 2, using  $\tau = 0.5$ . As illustrated in Tab. 4 and Tab. 5, the AGF method stochastically dominates the performance of these benchmark methods across all datasets, with the exception of the Yelp dataset.

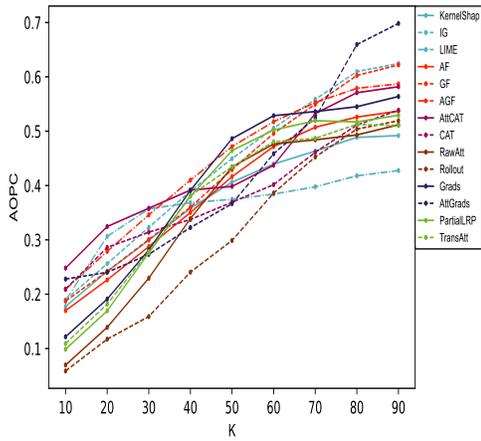
1034  
 1035  
 1036  
 1037  
 1038  
 1039  
 1040  
 1041  
 1042  
 1043  
 1044  
 1045  
 1046  
 1047  
 1048  
 1049  
 1050  
 1051  
 1052  
 1053  
 1054  
 1055  
 1056  
 1057  
 1058  
 1059  
 1060  
 1061  
 1062  
 1063  
 1064  
 1065  
 1066  
 1067  
 1068  
 1069  
 1070  
 1071  
 1072  
 1073  
 1074  
 1075  
 1076  
 1077  
 1078  
 1079  
 1080



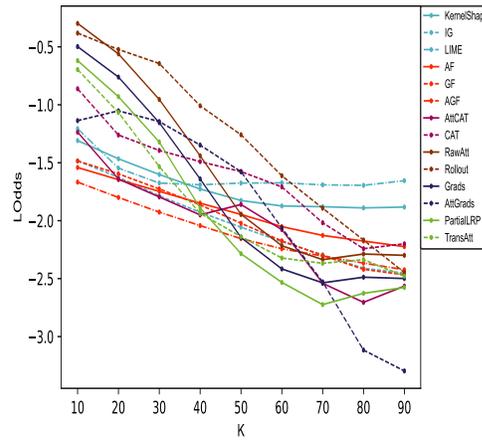
(a) AOPC score for IMDB dataset



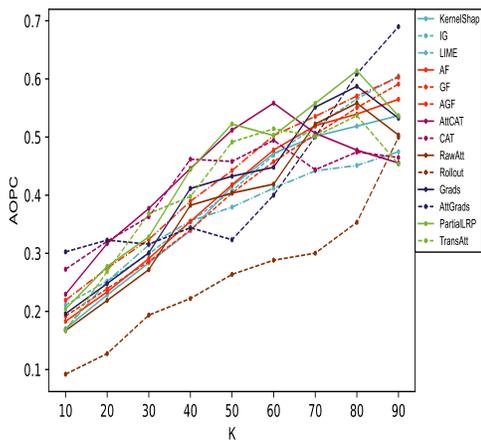
(b) LOdds score for IMDB dataset



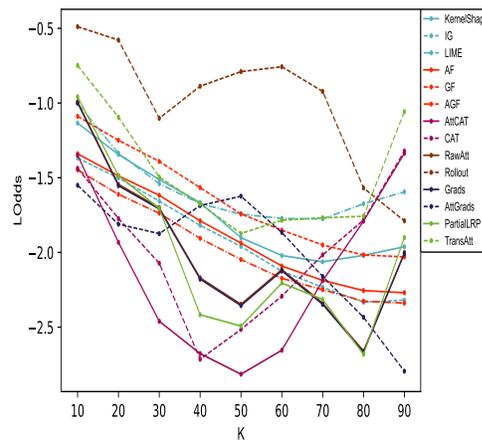
(c) AOPC score for Amazon dataset



(d) LOdds score for Amazon dataset



(e) AOPC score for Yelp dataset



(f) LOdds score for Yelp dataset

Figure 8: AOPC and LOdds scores of different methods in explaining BERT across the varying corruption rates  $k$  on IMDB, Amazon, and Yelp datasets. The x-axis illustrates masking the  $k\%$  of the tokens in an order of decreasing saliency.

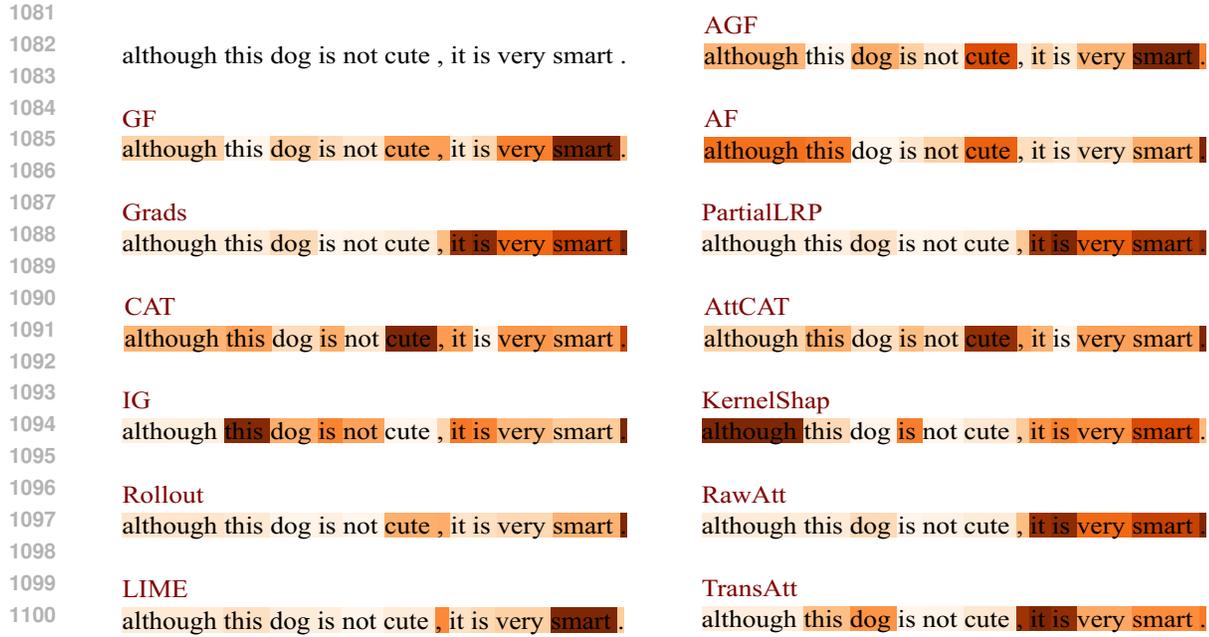


Figure 9: Visualizations of the normalized feature attributions generated by the methods on the binary classification task.

Table 4: ASO test to compare AGF method with the best benchmark method when we mask **top**  $k\%$  tokens.

Dataset	SST2		IMDB		Yelp		Amazon		AG News	
	AOPC	LOdds	AOPC	LOdds	AOPC	LOdds	AOPC	LOdds	AOPC	LOdds
$\epsilon_{\min}$	0.054	0.039	0.078	0.061	0.813	0.924	0.053	0.043	0.091	0.076

Table 5: ASO test to compare AGF method with the best benchmark method when we mask **bottom**  $k\%$  tokens.

Dataset	SST2		IMDB		Yelp		Amazon		AG News	
	AOPC	LOdds	AOPC	LOdds	AOPC	LOdds	AOPC	LOdds	AOPC	LOdds
$\epsilon_{\min}$	0.049	0.037	0.113	0.085	0.913	0.824	0.062	0.053	0.091	0.067

## E PROOFS

**Corollary 3.1.** While Shapley values are inherently unique, our findings in [Sec. 3.3](#) and [App. C](#) expose a critical inconsistency. We demonstrate that the optimal solution of the maximum flow problems defined in [eq. 8](#) is not necessarily unique, thereby disproving the claim that the feature attributions proposed by [Ethayarajh & Jurafsky \(2021\)](#) are Shapley values.

The non-uniqueness of these attributions, as evidenced by our proof, fundamentally conflicts with the defining properties of Shapley values. If these attributions defined by [Ethayarajh & Jurafsky \(2021\)](#) were indeed Shapley values, they would necessarily be unique. However, our observations demonstrate that since the optimal solution of the maximum flow problem is not necessarily unique, for each optimal solution of the maximum flow problem, we can derive corresponding feature attributions that differ from one another.

1128  
 1129  
 1130  
 1131  
 1132  
 1133  
 1134  
 1135  
 1136  
 1137  
 1138  
 1139  
 1140  
 1141  
 1142  
 1143  
 1144  
 1145  
 1146  
 1147  
 1148  
 1149  
 1150  
 1151  
 1152  
 1153  
 1154  
 1155  
 1156  
 1157  
 1158  
 1159  
 1160  
 1161  
 1162  
 1163  
 1164  
 1165  
 1166  
 1167  
 1168  
 1169  
 1170  
 1171  
 1172  
 1173  
 1174

□

**Theorem 3.1.** Since the optimal flow  $f^*$  is computed once for the entire graph and not for each potential subgraph, and the players (tokens) are all disjoint without any connections in  $S$ , blocking the flow through one player does not impact the outflow of any other players. Therefore, for every  $S \subseteq N$  where  $i \notin S$ , we have  $|f_{\text{out}}(i)| = v(S \cup \{i\}) - v(S)$ . Utilizing the definition of Shapley values in eq. 14, we obtain:

$$\begin{aligned}
 \phi_i(\vartheta) &= \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(|N| - |S| - 1)!}{|N|!} (\vartheta(S \cup \{i\}) - \vartheta(S)) \\
 &= \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(|N| - |S| - 1)!}{|N|!} (|f_{\text{out}}(i)|) \\
 &= |f_{\text{out}}(i)| \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(|N| - |S| - 1)!}{|N|!} \\
 &= |f_{\text{out}}(i)|
 \end{aligned}$$

It is also evident that the defined function meets all four fairness-based axioms of efficiency, symmetry, linearity, and additivity. □

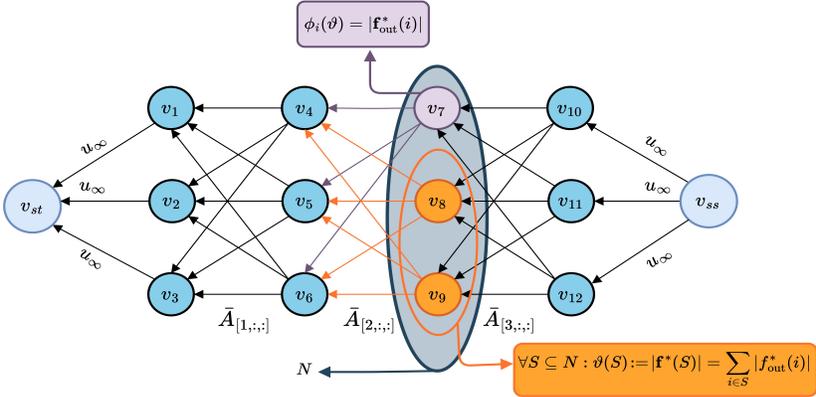


Figure 10: Visualizations of the procedure to define the cooperative game  $(N, \vartheta)$  using the solution of barrier-regularized maximum flow or its corresponding MCC problem.

## F IMPLEMENTATION DETAILS

### F.1 DATASETS

Tab. 6 illustrates comprehensive statistics of the datasets utilized for the classification task. We randomly extracted 5,000 sentences from each test section of the datasets, except for those with a test size less than 5,000, where we retained all samples. Furthermore, we prioritized diversity in our sampling process by incorporating sentences of varying lengths, with an equal distribution between those shorter and longer than the mode size of the test dataset.

Table 6: Statistical information and the pre-trained models employed for each dataset.

Datasets	# Test Samples	# Classes	$\ell_{\text{mode}}$	$\ell_{\text{min}}$	$\ell_{\text{max}}$	$\ell_{\text{avg}}$	Pre-trained Model
SST2	1,821	2	108	5	256	103.3	textattack/bert-base-uncased-SST-2
Amazon	5,000	2	127	15	1009	404.9	fabriceyhc/bert-base-uncased-amazon_polarity
IMDB	5,000	2	670	32	12988	1293.8	fabriceyhc/bert-base-uncased-imdb
Yelp	5,000	2	313	4	5107	723.8	fabriceyhc/bert-base-uncased-yelp_polarity
AG News	5,000	4	238	100	892	235.3	fabriceyhc/bert-base-uncased-ag_news

## F.2 TIME COMPLEXITY OF PROPOSED METHODS

In the minimum-cost circulation problem, we are given a directed graph  $G = (V, E)$  with  $|V| = n$  vertices and  $|E| = m$  edges, upper and lower edge capacities  $\mathbf{u}, \mathbf{l} \in \mathbb{R}^m$ , and edge costs  $\mathbf{c} \in \mathbb{R}^m$ . Our objective is to find a circulation  $\mathbf{f} \in \mathbb{R}^m$  satisfying:

$$\begin{aligned} \arg \min \mathbf{c}^\top \mathbf{f} \\ \mathbf{B}^\top \mathbf{f} = \mathbf{0} \\ \mathbf{l} \leq \mathbf{f} \leq \mathbf{u} \end{aligned}$$

where  $\mathbf{B} \in \mathbb{R}^{m \times n}$  is the edge-vertex incidence matrix.

To compare running times, we assume that  $\tilde{\mathbf{l}}, \tilde{\mathbf{u}}, \tilde{\mathbf{c}}$  represent the integral versions of  $\mathbf{l}, \mathbf{u}, \mathbf{c}$  obtained from either Algorithm 1 or Algorithm 2, and define  $U = \|\tilde{\mathbf{u}}\|_\infty$  and  $C = \|\tilde{\mathbf{c}}\|_\infty$ . Tab. 7 compares the latest iterative algorithms for solving the maximum flow and minimum-cost circulation problems. While the most efficient algorithm achieves nearly linear running time relative to the number of edges  $m$ , the runtime can still be significant with long input sequences. To solve the minimum-cost circulation problem, we implemented the algorithm by (Lee & Sidford, 2014) using CVXPY (Agrawal et al., 2018; Diamond & Boyd, 2016).

Table 7: Overview of recent iterative algorithms for maximum flow and minimum-cost circulation problems.

Year	MCC Bound	Max-Flow Bound	Author
2014	$O(m\sqrt{n} \text{polylog}(n) \log^2(U))$	$O(m\sqrt{n} \text{polylog}(n) \log^2(U))$	Lee & Sidford (2014)
2022	$O\left(m^{\frac{3}{2}} - \frac{1}{762} \text{polylog}(n) \log(U + C)\right)$	$O\left(m^{\frac{10}{7}} \text{polylog}(n) U^{\frac{1}{7}}\right)$	Axiotis et al. (2022)
2023	$O\left(m^{\frac{3}{2}} - \frac{1}{58} \text{polylog}(n) \log^2(U)\right)$	$O\left(m^{\frac{3}{2}} - \frac{1}{58} \text{polylog}(n) \log^2(U)\right)$	van den Brand et al. (2023)
2023	$O\left(m^{1+o(1)} \log(U) \log(C)\right)$	$O\left(m^{1+o(1)} \log(U) \log(C)\right)$	Chen et al. (2023a)

Experiments were conducted on a computing device running Ubuntu 20.04.4 LTS, equipped with an Intel(R) Xeon(R) Platinum 8368 CPU at 2.40GHz, featuring 12 cores and 24 threads for parallel processing. Graphics processing was handled by an NVIDIA RTX 3090 Ti with 40GB of dedicated memory. The device also had 230GB of system memory, ensuring ample computational resources for efficient and effective experimentation.

Tab. 8 compares the runtime of benchmark methods for analyzing the sentence "although this dog is not cute, it is very smart." Methods like RawAtt and Rollout, which rely on raw attention weights, have the shortest runtime. In contrast, methods requiring complex post-processing to evaluate feature attributions have longer runtime. Our proposed methods' runtime is comparable to others in this latter group.

Table 8: Runtime of all methods for the showcase example.

Methods	Runtime (seconds)
RawAtt	0.123
Rollout	0.154
Grads	1.554
AttGrads	1.571
CAT	1.684
AttCAT	1.660
PartialLRP	1.571
TransAtt	1.620
LIME	1.462
KernelShap	2.342
IG	2.701
AF	2.301
GF	2.305
AGF	2.306