

---

# UNIFYING FEATURE ATTRIBUTION: BRIDGING ATTENTION MECHANISMS AND SHAPLEY VALUES

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Machine learning (ML) applications have become ubiquitous in various domains, including and large language models and computer vision. However, the increasing prevalence of ML models has also emphasized the need for explainable artificial intelligence (XAI). In particular, it has become crucial to develop techniques that can help explain how these models work and why they make certain decisions. This is especially important in cases where these models are used to make critical decisions that can have significant consequences on individuals or society as a whole.

## 1 INTRODUCTION

In the pursuit of transparent and interpretable machine learning models, understanding feature attribution methods is pivotal. Feature attribution elucidates the contribution of individual input elements to model predictions, fostering model interpretability. However, existing methods grapple with challenges, such as non-uniqueness and interpretability gaps. In this paper, we present a cohesive framework by leveraging attention mechanisms and Shapley values.

The attention mechanism, initially designed for natural language processing, has proven instrumental in capturing contextual relationships within input data. Simultaneously, Shapley values, rooted in cooperative game theory, provide a principled approach to distributing contributions among features. We propose a novel integration of attention and Shapley values, unifying their strengths to yield robust and interpretable feature attributions.

### **Contribution:**

Our approach addresses the challenges faced by existing methods, offering a more nuanced understanding of feature importance. Through empirical evaluations on diverse datasets, we demonstrate the effectiveness and versatility of our proposed method in enhancing interpretability across various machine learning models.

## 2 RELATED WORKS

The concept of attention flows was introduced by Abnar & Zuidema (2020) and presented a new perspective on how XAI could leverage attention flows. However, the authors of the paper chose to focus on simple ideas that only required attention weights, which could be easily implemented in any architecture employing self-attention mechanism. While they compared attention flows with other concepts like attention rollout, they did not employ maximum flow to obtain the attribution of each token in language models. In the classification task, they conducted an evaluation of the weights obtained from raw attention, attention rollout, and attention flow for the CLS embedding over input tokens in all attention layers for a selected set of examples. In the masked language modeling task, a similar analysis was performed on the quantities derived from the MASK embedding to the two potential references.

In a recent study, Ethayarajh & Jurafsky (2021) attempted to establish a systematic connection between attention flows and XAI. While their work focused on demonstrating that attention flows can be interpreted as Shapley values under certain conditions, they did not account for the non-uniqueness of such flows. It is worth noting that, although maximum flow problem has the unique optimal value, there can be multiple feasible flows that attain this optimal value.

This highlights the need for further research to explore the method to resolve the issue of non-uniqueness in attention flows and its impact on XAI.

### 3 PRELIMINARIES

We will define preliminary terms and definitions here that will be used in the subsequent sections. This study will primarily utilize these concepts.

#### 3.1 GRAPH NETWORK FLOW

##### 3.1.1 MAXIMUM FLOW

Given a network  $G = (V, E, \mathbf{u})$ , where  $u(i, j)$  is the capacity for the edge  $(i, j) \in E$ , a flow is characterized as a function  $f : E \rightarrow \mathbb{R}^{\geq 0}$  s.t.

$$\begin{aligned} f_{ij} &\leq u_{ij} \quad \forall (i, j) \in E \quad (\text{capacity constraints}) \\ \sum_{j:(i,j) \in E} f_{ij} - \sum_{j:(j,i) \in E} f_{ji} &= 0, \quad \forall i \in V, i \neq s, t \quad (\text{flow conservation constraints}) \end{aligned} \quad (1)$$

The value of a flow in a given network  $G = (V, E, s, t, \mathbf{u})$  is established as the sum of flows into the source vertex  $s$ , denoted as  $\text{val}(f) = \sum_{v:(v,s)} f_{vs}$ , and a maximum flow is identified as a feasible flow with the highest attainable value.

Behrooz: Check definition of val(f)

##### 3.1.2 MULTI-COMMODITY MAXIMUM FLOW

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

$$\sum_{k=1}^r f^k(i, j) \leq u(i, j) \quad \forall (i, j) \in E \quad (2)$$

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

To solve the problem of multi-commodity maximum flow, we can simplify it by transforming it into a standard maximum flow problem. This can be achieved by introducing two new nodes - a "super-source" node  $s$  and a "super-sink" node  $t$ . The "super-source" node  $s$  should be connected to all the original sources

$s_i$  through edges of finite capacities, while the "super-sink" node  $t$  should be connected to all the original sinks  $t_i$  with edges of finite capacities:

- Each outgoing edge from the "super-source" node  $s$  to each source node  $s_i$  gets assigned a capacity that is equal to the total capacity of the outgoing edges from the source node  $s_i$ .
- Each incoming edge from an original "super-sink" node  $t$  to each sink node  $t_i$  gets assigned a capacity that is equal to the total capacity of the incoming edge to the sink node  $t_i$ .

It is easy to see that the maximum flow from  $s$  to  $t$  is equivalent to the maximum sum of flows in a feasible multi-commodity flow in the original network.

### 3.1.3 MINIMUM-COST CIRCULATION

Given a network  $G = (V, E, \mathbf{u}, \mathbf{l}, \mathbf{c})$ , where  $c_{ij}$  is the cost,  $l_{i,j}$  and  $u_{i,j}$  are the capacities for the edge  $(i, j) \in E$ , a circulation is a function  $f : E \rightarrow \mathbb{R}^{\geq 0}$  s.t.

$$\begin{aligned} l_{ij} &\leq f_{ij} \leq u_{ij}, & \forall (i, j) \in E \\ \sum_{j:(i,j) \in E} f_{ij} - \sum_{j:(j,i) \in E} f_{ji} &= 0, & \forall i \in V. \end{aligned} \quad (3)$$

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

Given a graph  $G = (V, E)$  with  $|V| = n$  vertices and  $|E| = m$  edges, the edge-vertex incidence matrix is denoted as  $B_{m \times n}$ . For a directed graph, the entries of the matrix  $B$  are defined as follows:

$$b_{ev} = \begin{cases} -1, & \text{if vertex } v \text{ is the tail of edge } e, \\ 1, & \text{if vertex } v \text{ is the head of edge } e, \\ 0, & \text{if edge } e \text{ is not incident to vertex } v. \end{cases} \quad (4)$$

Then, the formulation of minimum-cost circulation problem can be written as the following linear programming (LP) problem:

$$\begin{aligned} \min_{\substack{\mathbf{B}^\top \mathbf{f} = \mathbf{0} \\ l_e \leq f_e \leq u_e \forall e \in E}} \mathbf{c}^\top \mathbf{f} \quad \text{i.e.} \quad \arg \min_{\substack{\mathbf{B}^\top \mathbf{f} = \mathbf{0} \\ l \leq \mathbf{f} \leq \mathbf{u}}} \mathbf{c}^\top \mathbf{f} \end{aligned} \quad (5)$$

### 3.1.4 SHAPLEY VALUES

Shapley values are defined through value functions  $\vartheta$  of subsets  $S$  of players in a set  $N = \{1, 2, \dots, n\}$ . Specifically, the Shapley value  $\phi_j(\vartheta)$  of player (or feature)  $j$  in a game with total payoff  $v$  is given by:

$$\phi_j(\vartheta) = \sum_{S \subseteq N \setminus \{j\}} \frac{|S|!(n - |S| - 1)!}{n!} (\vartheta(S \cup \{j\}) - \vartheta(S)) \quad (6)$$

These Shapley values are the only constructs that jointly satisfy the Efficiency, Symmetric, Null Player and Additivity axioms:

- 
- **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  $\phi_i(\vartheta) = \phi_j(\vartheta)$ .
  - **Null Player:** 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  $\phi_i(\vartheta) = 0$ .
  - **Additivity:** If two games are combined into a larger game, the Shapley value of a player in the combined game is the sum of their Shapley values in the individual games.

### 3.1.5 ATTENTION FLOWS

Incorporating attention mechanisms into XAI can be accomplished using a flow network representation of the attention graph. In this representation, the graph's edges are given capacities corresponding to attention weights. By applying a maximum flow algorithm, we can determine the maximum attention flow from any node in any layer to any input node, which can be used to estimate the amount of attention directed towards the input nodes.

However, to evaluate the maximum flow from all output nodes to all input nodes, we need to implement the concept of multi-commodity flow. This involves introducing a super-source  $s$  and a super-target  $t$  with a  $u_\infty$  capacity. The interconnection between layers is established using attention weights, creating a layered network.

To formalize this process, we create the attention flow matrix  $\mathcal{F}$  and its integral version  $\tilde{\mathcal{F}}$  by mapping the network, which consists of  $T$  tokens and  $L$  layers, according to the Algorithms 1 and 2:

---

#### Algorithm 1 Get Attention Flow-Backward Approach

---

**Require:**  $\mathcal{A}_{L \times T \times T}$ : A tensor containing the aggregated attentions.

**Ensure:** Tuple:  $(\mathcal{F}, \tilde{\mathcal{F}})$

```

function GET_ATTENTION_FLOW( $\mathcal{A}$ )
   $L, T, _ \leftarrow \mathcal{A}.shape()$ 
   $\beta_{\min} \leftarrow \min(\mathcal{A} > 0)$ 
   $\beta \leftarrow \lfloor \log_{10}(\beta_{\min}) \rfloor$ 
   $\gamma \leftarrow 10^\beta$ 
   $\mathcal{F} \leftarrow \text{np.zeros}((L + 1) * T + 2, (L + 1) * T + 2)$ 
   $u_\infty \leftarrow T$ 
  for  $i$  in  $\text{range}(T)$  do
     $\mathcal{F}[i + 1][0] \leftarrow u_\infty$ 
  end for
  for  $i$  in  $\text{range}(T)$  do
     $\mathcal{F}[-1][-i - 2] \leftarrow u_\infty$ 
  end for
  for  $\ell$  in  $\text{range}(0, L)$  do
     $\mathcal{F}[T * (\ell + 1) + 1 : T * (\ell + 2) + 1, T * (\ell) + 1 : T * (\ell + 1) + 1] \leftarrow \mathcal{A}[\ell]$ 
  end for
   $\tilde{\mathcal{F}} \leftarrow \text{int}(\gamma * \mathcal{F})$ 
end function

```

---

---

**Algorithm 2** Get Attention Flow-Forward Approach

---

**Require:**  $\mathcal{A}_{L \times T \times T}$ : A tensor containing the aggregated attentions.

**Ensure:** Tuple:  $(\mathcal{F}, \tilde{\mathcal{F}})$

```
function GET_ATTENTION_FLOW( $\mathcal{A}$ )  
   $L, T, _ \leftarrow \mathcal{A}.\text{shape}()$   
   $\beta_{\min} \leftarrow \min(\mathcal{A} > 0)$   
   $\beta \leftarrow -\lfloor \log_{10}(\beta_{\min}) \rfloor$   
   $\gamma \leftarrow 10^\beta$   
   $\mathcal{F} \leftarrow \text{np.zeros}((L + 1) * T + 2, (L + 1) * T + 2)$   
   $u_\infty \leftarrow T$   
  for  $i$  in  $\text{range}(T)$  do  
     $\mathcal{F}[i + 1][0] \leftarrow u_\infty$   
  end for  
  for  $i$  in  $\text{range}(T)$  do  
     $\mathcal{F}[-1][-i - 2] \leftarrow u_\infty$   
  end for  
  for  $\ell$  in  $\text{range}(0, L)$  do  
     $\mathcal{F}[T * (\ell + 1) + 1 : T * (\ell + 2) + 1, T * (\ell) + 1 : T * (\ell + 1) + 1] \leftarrow \mathcal{A}[\ell]$   
  end for  
   $\tilde{\mathcal{F}} \leftarrow \text{int}(\gamma * \mathcal{F})$   
end function
```

---

Behrooz: check the algorithms accuracy add source  
and node indices to algorithms

## 4 METHODS

### 4.1 NON-UNIQUENESS OF MAXIMUM FLOW

It should be noted that the problem at hand is not strictly convex, which implies that the existence of a unique solution cannot be guaranteed. Additionally, we have discovered that the maximum flow problem for the previously constructed graph does not possess a unique solution either. Specifically, we have identified two feasible maximum flows for a similar graph created using attention flows. These maximum flows can be obtained using two distinct algorithms, namely Algorithms 1 and 2.

Moreover, we know that if any linear programming has two distinct feasible solutions  $f_1, f_2$ , any convex combination  $\gamma_1 f_1 + \gamma_2 f_2$  where  $\gamma_1 + \gamma_2 = 1$  is a new solution. Consequently, the attention flows has infinite feasible solution and therefore the feature attributions and Shapley values resulted from this approach is ill-defined.

## 5 MAKE MAXIMUM FLOW UNIQUE

As previously discussed, each network inherently possesses a distinct maximum flow value. However, achieving this maximum value is not confined to a singular feasible flow; rather, multiple feasible flows can attain the same maximum. This variability introduces a challenge in discerning the optimal flow solution. To address this challenge, the introduction of regularization terms into the maximum flow problem proves beneficial. Through this regularization, the objective of the linear programming problem defined in Eq. 9 is transformed into a strictly convex function. This transformation ensures the derivation of a unique and optimal approximate solution, thereby overcoming the issue posed by the non-uniqueness of feasible flows achieving the maximum value.

Here we introduce two main regularization term and the logic behind them: 1. log-barrier regularization 2.  $\ell_1$  and  $\ell_2$  regularization.

---

## 5.1 LOG-BARRIER REGULARIZATION

We can reformulate the minimum cost circulation problem as:

$$\begin{aligned} \arg \min \quad & \mathbf{c}^\top \mathbf{f} \\ \text{subject to} \quad & \mathbf{B}^\top \mathbf{f} = \mathbf{0} \\ & h(\mathbf{f}) \leq \mathbf{0} \end{aligned} \tag{7}$$

where  $h(\mathbf{f}) = (\mathbf{f} - \mathbf{l})(\mathbf{f} - \mathbf{u})$ .

Defining the log-barrier function  $\phi_\mu(\mathbf{f}) = -\mu \sum_{e \in E} \log(-h(f_e)) = \mu \sum_{e \in E} (\log(f_e - l_e) + \log(u_e - f_e))$ , the original problem can, therefor, be approximated using the log barrier function as the following:

$$\begin{aligned} \arg \min \quad & \mathbf{c}^\top \mathbf{f} + \phi_\mu(\mathbf{f}) \\ \text{subject to} \quad & \mathbf{B}^\top \mathbf{f} = \mathbf{0} \end{aligned} \tag{8}$$

It is easy to observe that as long as  $\mu > 0$  and our initial solution is feasible, the barrier function will ensure that any solution we find using an iterative minimization scheme, such as gradient descent or Newton's method, remains feasible. Moreover, one can show that to get an  $\varepsilon$ -approximate solution to the Eq. 7 it suffices to set  $\mu \leq \frac{\varepsilon}{2m}$  and find the optimal solution to the corresponding problem in Eq. 8.

Finally, observe that the Hessian of the objective function of Eq. 8 at some point  $\mathbf{f}$  is equal to the Hessian  $\nabla^2 \phi_\mu(x)$  of the barrier function. The latter turns out to be a diagonal matrix with

$$(\nabla^2 \phi_\mu(\mathbf{f}))_{e,e} = \mu \left( \frac{1}{(u_e - f_e)^2} + \frac{1}{(f_e - l_e)^2} \right),$$

for each edge  $e$ . Thus, the Hessian of our objective function is positive definite (provided  $\mu > 0$ ), which implies that the objective function is strongly convex and therefore have a unique feasible solution.

There are many preceding study about various numerical methods and their convergence rate which we leave the readers for further investigation (Bertsimas & Tsitsiklis, 1997; Boyd & Vandenberghe, 2004; Bubeck, 2015). It is noteworthy to mention that many other barrier functions can be used instead of log-barrier functions which can enhance the convergence rate of numerical solver for the optimization problem. (Cohen et al., 2020; Lee & Sidford, 2020; Axiotis et al., 2022; van den Brand et al., 2021; Chen et al., 2022)

### 5.1.1 $\ell_1$ AND $\ell_2$ REGULARIZATION

The similar idea is to apply  $\ell_1$  and  $\ell_2$  regularization.

---

## REFERENCES

- Samira Abnar and Willem Zuidema. Quantifying Attention Flow in Transformers, May 2020.
- 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.
- Dimitris Bertsimas and John N. Tsitsiklis. *Introduction to Linear Optimization*. Athena Scientific Series in Optimization and Neural Computation. Athena Scientific, Belmont, Mass, 1997. ISBN 978-1-886529-19-9.
- 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.
- Li Chen, Rasmus Kyng, Yang P. Liu, Richard Peng, Maximilian Probst Gutenberg, and Sushant Sachdeva. Maximum Flow and Minimum-Cost Flow in Almost-Linear Time, April 2022.
- Michael B. Cohen, Yin Tat Lee, and Zhao Song. Solving Linear Programs in the Current Matrix Multiplication Time, October 2020.
- Kawin Ethayarajh and Dan Jurafsky. Attention Flows are Shapley Value Explanations, May 2021.
- Yin Tat Lee and Aaron Sidford. Solving Linear Programs with  $\sqrt{\text{rank}}$  Linear System Solves, August 2020.
- Jan van den Brand, Yin Tat Lee, Yang P. Liu, Thatchaphol Saranurak, Aaron Sidford, Zhao Song, and Di Wang. Minimum Cost Flows, MDPs, and  $\ell_1$ -Regression in Nearly Linear Time for Dense Instances, August 2021.

## A APPENDIX

You may include other additional sections here.

### A.1 DEFINITIONS

#### A.1.1 MINIMUM COST CIRCULATION PROBLEM

the formulation of minimum-cost circulation problem can be written as the following linear programming (LP) problem:

$$\min_{\substack{\mathbf{B}^\top \mathbf{f} = \mathbf{0} \\ l_e \leq f_e \leq u_e \forall e \in E}} \mathbf{c}^\top \mathbf{f} \quad \text{i.e.} \quad \arg \min_{\substack{\mathbf{B}^\top \mathbf{f} = \mathbf{0} \\ \mathbf{l} \leq \mathbf{f} \leq \mathbf{u}}} \mathbf{c}^\top \mathbf{f} \quad (9)$$

It is easy to show that it satisfies strong duality with dual problem:

$$\max_{\substack{\mathbf{B}\mathbf{y} + \mathbf{s}_l - \mathbf{s}_u = \mathbf{c} \\ \mathbf{s}_l, \mathbf{s}_u \in \mathbb{R}_{\geq 0}^E}} -\mathbf{s}_u^\top \mathbf{u} + \mathbf{s}_l^\top \mathbf{l} \quad (10)$$

where the Lagrangian ( $\mathcal{L}$ ) is defined as:

$$\mathcal{L}(\mathbf{f}, \mathbf{y}, \mathbf{s}_u, \mathbf{s}_l) = \mathbf{c}^\top \mathbf{f} + \mathbf{y}^\top (\mathbf{B}^\top \mathbf{f}) + \mathbf{s}_u^\top (\mathbf{u} - \mathbf{f}) + \mathbf{s}_l^\top (\mathbf{f} - \mathbf{l})$$

Since our primal LP problem is convex, the KKT conditions provide both necessary and sufficient conditions to find optimal primal and dual variables for the minimum cost circulation problem as follows:

1. Stationarity:

$$\nabla \mathcal{L}(\mathbf{f}, \mathbf{y}, \mathbf{s}_u, \mathbf{s}_l) = \mathbf{c} - \mathbf{B}\mathbf{y} + \mathbf{s}_u - \mathbf{s}_l = \mathbf{0}$$

2. Complementary Slackness:

$$\mathbf{s}_u^\top (\mathbf{u} - \mathbf{f}) = \mathbf{0} : s_e^+ \cdot (u_e - f_e) = 0, \quad \forall e \in E \quad (\text{for upper-bound capacity constraints})$$

$$\mathbf{s}_l^\top (\mathbf{f} - \mathbf{l}) = \mathbf{0} : s_e^- \cdot (f_e - l_e) = 0, \quad \forall e \in E \quad (\text{for lower-bound capacity constraints})$$

3. Primal Feasibility:

$$\mathbf{B}^\top \mathbf{f} = \mathbf{0} \quad (\text{flow conservation constraint})$$

$$\mathbf{l} \leq \mathbf{f} \leq \mathbf{u} : l_e \leq f_e \leq u_e, \quad \forall e \in E \quad (\text{capacity constraints})$$

4. Dual Feasibility:

$$\mathbf{s}_u, \mathbf{s}_l \geq \mathbf{0} \quad (\text{non-negativity of slack variables})$$

Here,  $\mathbf{y}$  represents the dual variables corresponding to the equality constraint ( $\mathbf{B}^\top \mathbf{f} = \mathbf{0}$ ), and  $\mathbf{s}_u$  and  $\mathbf{s}_l$  represent the dual variables corresponding to the upper and lower inequality constraints ( $l_e \leq f_e \leq u_e$ ), respectively.