

GlobEnc: Quantifying Global Token Attribution by Incorporating the Whole Encoder Layer in Transformers

Anonymous ACL submission

Abstract

There has been a growing interest in interpreting the underlying dynamics of Transformers. While self-attention patterns were initially deemed as the primary choice, recent studies have shown that integrating other components can yield more accurate explanations. This paper introduces a novel token attribution analysis method that incorporates all the components in the encoder block and aggregates this throughout layers. We quantitatively and qualitatively demonstrate that our method can yield faithful and meaningful global token attributions. Our extensive experiments reveal that incorporating almost every encoder component results in increasingly more accurate analysis in both local (single layer) and global (the whole model) settings. Our global attribution analysis surpasses previous methods by achieving significantly higher results in various datasets.

1 Introduction

The stellar performance of Transformers (Vaswani et al., 2017) has garnered a lot of attention to analyzing the reasons behind their effectiveness. The self-attention mechanism has been one of the main areas of focus (Clark et al., 2019; Kovaleva et al., 2019; Reif et al., 2019; Htut et al., 2019). But, there have been debates on whether raw attention weights are reliable anchors for explaining model’s behavior (Wiegrefe and Pinter, 2019; Serrano and Smith, 2019; Jain and Wallace, 2019). Recently, it was shown that incorporating vector norms should be an indispensable part of any attention-based analysis (Kobayashi et al., 2020, 2021). However, these norm-based studies incorporate only the attention block into their analysis, whereas the encoder layer is composed of more components. We show that these components are essential for a more accurate analysis. Moreover, these studies are constrained to the analysis of single layer attributions.

In order to expand the analysis to the entire model, an aggregation technique has to be em-

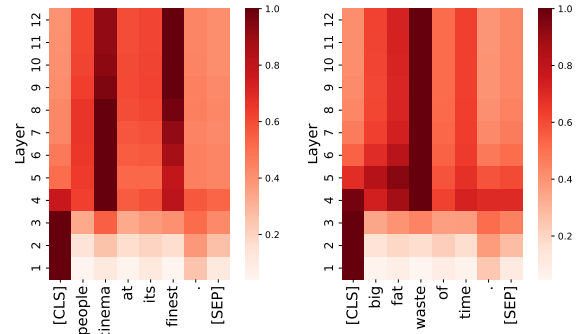


Figure 1: Aggregated attribution maps (\mathcal{N}_{ENC}) for the [CLS] token for fine-tuned BERT on SST2 dataset (sentiment analysis). Our method (GlobEnc) is able to accurately quantify the global attribution of the model.

ployed. Abnar and Zuidema (2020) proposed two aggregation methods, *rollout* and *max-flow*, which combine raw attention weights across layers. Despite reporting improvements on the attributions, the final results are still very low on fine-tuned models. Also, gradient-based alternatives have been argued to provide a more robust basis for such analysis (Brunner et al., 2020; Pascual et al., 2021), while being computationally intensive.

In this paper, we propose a new global token attribution analysis method (GlobEnc), which incorporates not only the attention block, but also the second layer normalization that produces the encoder layer’s output. Our results on BERT (Devlin et al., 2019) show high correlations with gradient based methods in both local and global settings.

To evaluate our approach, we compare the global attribution with the input token attributions obtained by gradient-based saliency scores. We show that: (i) norm-based methods achieve higher correlation than weight-based methods; (ii) incorporating residual connections plays an essential role in token attribution; (iii) layer normalizations can improve our analysis only if coupled together; and (iv) aggregation across layers is crucial for an ac-

curate whole-model attribution analysis. Based on these findings, we propose a global attribution method that provides faithful and plausible results (Figure 1). In summary, our main contributions are threefold:

- We expand the scope of analysis from attention block in Transformers to the whole encoder.
- Our method significantly improves over existing techniques for quantifying global token attribution in BERT.
- We qualitatively demonstrate that the attributions obtained by our method are plausibly interpretable.

2 Background

In encoder-based language models (such as BERT), a Transformer encoder layer is composed of several components (Figure 2). The core component of the encoder is the self-attention mechanism (Vaswani et al., 2017), which is responsible for the information mixture of a sequence of token representations ($\mathbf{x}_1, \dots, \mathbf{x}_n$). Each self-attention head computes a set of attention weights $\mathbf{A}^h = \{\alpha_{i,j}^h | 1 \leq i, j \leq n\}$, where $\alpha_{i,j}^h$ is the raw attention weight from the i^{th} token to the j^{th} token in head $h \in \{1, \dots, H\}$. Therefore, the output representation ($\mathbf{z}_i \in \mathbb{R}^d$) for the i^{th} token of a multi-head (with H heads) self-attention module is computed by concatenating the heads' outputs followed by a head-mixing \mathbf{W}_O projection:

$$\mathbf{z}_i = \text{CONCAT}(\mathbf{z}_i^1, \dots, \mathbf{z}_i^H) \mathbf{W}_O \quad (1)$$

Where each head's output vector is generated by performing a weighted sum over the transformed value vectors $\mathbf{v}(\mathbf{x}_j) \in \mathbb{R}^{d_v}$:

$$\mathbf{z}_i^h = \sum_{j=1}^n \alpha_{i,j}^h \mathbf{v}^h(\mathbf{x}_j) \quad (2)$$

Norm-based attention. While one may interpret the attention mechanism using the attention weights \mathbf{A} , Kobayashi et al. (2020) argued that doing so would ignore the norm of the transformed vectors multiplied by the weights, elucidating that the weights are insufficient for interpretation. Their solution enhanced the interpretability of attention weights by incorporating the value vectors $\mathbf{v}(\mathbf{x}_j)$ and the following projection \mathbf{W}_O . By reformulating Equation 1, we can consider \mathbf{z}_i as a summation

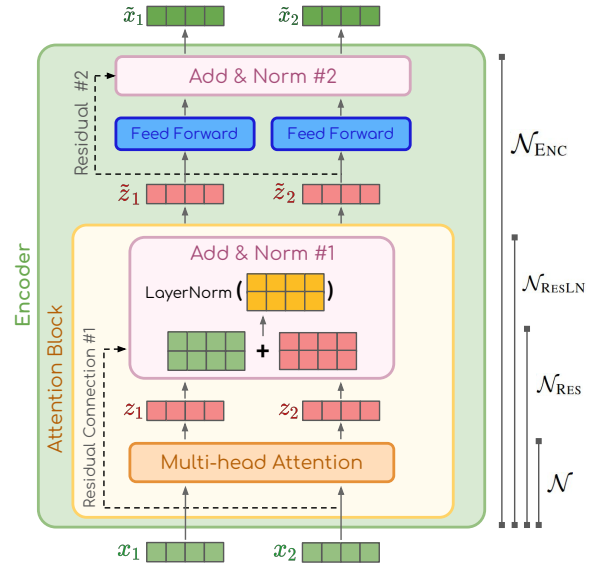


Figure 2: The internal structure of a Transformer encoder layer. We show on the diagram the components that are incorporated by each token attribution analysis method. Our method incorporates the whole encoder (\mathcal{N}_{ENC}) except for the fully connected feed-forward module. Diagram inspired by Alammari (2018).

over the attentions heads:

$$\mathbf{z}_i = \sum_{h=1}^H \sum_{j=1}^n \alpha_{i,j}^h \underbrace{\mathbf{v}^h(\mathbf{x}_j) \mathbf{W}_O^h}_{f^h(\mathbf{x}_j)} \quad (3)$$

Using this reformulation¹, Kobayashi et al. proposed a *norm-based* token attribution analysis method, $\mathcal{N} := (\|\mathbf{z}_{i \leftarrow j}\|) \in \mathbb{R}^{n \times n}$, to measure each token's contribution in a self-attention module:

$$\mathbf{z}_{i \leftarrow j} = \sum_{h=1}^H \alpha_{i,j}^h f^h(\mathbf{x}_j) \quad (4)$$

They showed that incorporating the magnitude of the transformation function ($f^h(\mathbf{x})$) is crucial in assessing the input tokens' contribution to the self-attention output.

Residual connections & Layer Normalizations.

Kobayashi et al. (2021) added the attention block's Layer Normalization (LN #1) and Residual connection (RES #1) to its prior norm-based analysis to assess the impact of residual connections and layer normalization inside an attention block. $\mathcal{N}_{\text{RES}} := (\|\mathbf{z}_{i \leftarrow j}^+\|) \in \mathbb{R}^{n \times n}$ is the analysis method which incorporates the attention block's residual

¹ \mathbf{W}_O^h is a head-specific slice of the original \mathbf{W}_O projection. For more information about the reformulation process, see Appendix C in Kobayashi et al. (2021)

connection. The input vector \mathbf{x} is added to the attribution of each token to itself to incorporate the influence of residual connection #1:

$$\mathbf{z}_{i \leftarrow j}^+ = \sum_{h=1}^H \alpha_{i,j}^h f^h(\mathbf{x}_i) + \mathbf{1}[i=j]\mathbf{x}_i \quad (5)$$

They proposed a method for decomposing LN² into a summation of normalizations:

$$\begin{aligned} \text{LN}(\mathbf{z}_i^+) &= \sum_{j=1}^n g_{\mathbf{z}_i^+}(\mathbf{z}_{i \leftarrow j}^+) + \beta \\ g_{\mathbf{z}_i^+}(\mathbf{z}_{i \leftarrow j}^+) &:= \frac{\mathbf{z}_{i \leftarrow j}^+ - m(\mathbf{z}_{i \leftarrow j}^+)}{s(\mathbf{z}_i^+)} \odot \gamma \end{aligned} \quad (6)$$

where $m(\cdot)$ and $s(\cdot)$ are the element-wise mean and standard deviation of the input vector (cf. §A.1). The decomposition can be applied to the contribution vectors:

$$\tilde{\mathbf{z}}_{i \leftarrow j} = g_{\mathbf{z}_i^+} \left(\sum_{h=1}^H \alpha_{i,j}^h f^h(\mathbf{x}_i) + \mathbf{1}[i=j]\mathbf{x}_i \right) \quad (7)$$

Accordingly, we can compute the magnitude $\mathcal{N}_{\text{RESLN}} := (\|\tilde{\mathbf{z}}_{i \leftarrow j}\|) \in \mathbb{R}^{n \times n}$, which represents the amount of influence of an encoder layer’s input token j on its output token i . Based on this formulation, a context-mixing ratio could be defined as:

$$r_i = \frac{\|\sum_{j=1, j \neq i}^n \tilde{\mathbf{z}}_{i \leftarrow j}\|}{\|\sum_{j=1, j \neq i}^n \tilde{\mathbf{z}}_{i \leftarrow j}\| + \|\tilde{\mathbf{z}}_{i \leftarrow i}\|} \quad (8)$$

Experiments by Kobayashi et al. (2021) revealed considerably low r values which indicates the huge impact of the residual connections. In other words, the model tends to preserve token representations more than mixing them with each other.

3 Methodology

Our method for input token attribution analysis has a holistic view and takes into account almost every component within the encoder layer. To this end, we first extend the norm-based analysis of Kobayashi et al. (2021) by incorporating the encoder’s output layer normalization #2. We then apply an aggregation technique to combine the information flow throughout all layers.

² $\gamma \in \mathbb{R}^d$ and $\beta \in \mathbb{R}^d$ are the trainable weights of LN. Similar to Kobayashi et al. (2021) we ignore β .

Encoder layer output \neq Attention block output.

While the residual connection #1 and the layer normalization #1 from the attention block are included in the analysis of Kobayashi et al. (2021), the subsequent FFN, residual connection #2, and output LN #2 are ignored (see Fig. 2). Hence, $\mathcal{N}_{\text{RESLN}}$ might not be indicative of the entire encoder layer’s function. To address this issue, we additionally include the encoder layer components from the attention block outputs ($\tilde{\mathbf{z}}_i$) to the output representations ($\tilde{\mathbf{x}}_i$). The output of each encoder ($\tilde{\mathbf{x}}_i$) is computed as follows:

$$\begin{aligned} \tilde{\mathbf{z}}_i^+ &= \text{FFN}(\tilde{\mathbf{z}}_i) + \tilde{\mathbf{z}}_i \\ \tilde{\mathbf{x}}_i &= \text{LN}(\tilde{\mathbf{z}}_i^+) \end{aligned} \quad (9)$$

We apply the LN decomposition rule in Eq. 7 to separate the impacts of residual and FFN output:

$$\tilde{\mathbf{x}}_i = \sum_{j=1}^n \left(g_{\tilde{\mathbf{z}}_i^+}(\text{FFN}(\tilde{\mathbf{z}}_{i \leftarrow j})) + g_{\tilde{\mathbf{z}}_i^+}(\tilde{\mathbf{z}}_{i \leftarrow j}) \right) + \beta \quad (10)$$

Given that the activation function between the two fully connected layers in the FFN component is non-linear (Vaswani et al., 2017), a linear decomposition similar to Eq. 7 cannot be derived. As a result, we omit FFN’s influence on the contribution of each token and instead consider residual connection #2, approximating $\tilde{\mathbf{x}}_{i \leftarrow j}$ as $g_{\tilde{\mathbf{z}}_i^+}(\tilde{\mathbf{z}}_{i \leftarrow j})$. Nevertheless, it should be noted that the FFN still preserves some influence on this new setting due to the presence of $s(\tilde{\mathbf{z}}_i^+)$ in $g_{\tilde{\mathbf{z}}_i^+}(\tilde{\mathbf{z}}_{i \leftarrow j})$. Similar to Eq. 7, we can introduce a more inclusive layerwise analysis method $\mathcal{N}_{\text{ENC}} := (\|\tilde{\mathbf{x}}_{i \leftarrow j}\|) \in \mathbb{R}^{n \times n}$ from input token j to output token i using:

$$\tilde{\mathbf{x}}_{i \leftarrow j} \approx g_{\tilde{\mathbf{z}}_i^+}(\tilde{\mathbf{z}}_{i \leftarrow j}) = \frac{\tilde{\mathbf{z}}_{i \leftarrow j} - m(\tilde{\mathbf{z}}_{i \leftarrow j})}{s(\tilde{\mathbf{z}}_i^+)} \odot \gamma \quad (11)$$

Aggregating multi-layer attention. To create an aggregated attribution score, Abnar and Zuidema (2020) proposed describing the model’s attentions via modelling the information flow with a directed graph. They introduced *attention rollout* method, which linearly combines attention along all available paths in the pairwise attention graph. The attention rollout of layer ℓ w.r.t. the inputs is computed recursively as follows:

$$\tilde{\mathbf{A}}_\ell = \begin{cases} \hat{\mathbf{A}}_\ell \tilde{\mathbf{A}}_{\ell-1} & \ell > 1 \\ \hat{\mathbf{A}}_\ell & \ell = 1 \end{cases} \quad (12)$$

$$\hat{\mathbf{A}}_\ell = 0.5\bar{\mathbf{A}}_\ell + 0.5\mathbf{I} \quad (13)$$

\bar{A}_ℓ is the raw attention map averaged across all heads in layer ℓ . This method assumes equal contribution from the residual connection and multi-head attention (See Fig. 2). Hence, an identity matrix is summed and renormalized, giving \hat{A}_ℓ .

For aggregating the layerwise analysis methods, we use the rollout technique with minor modifications. As many of the methods already include residual connections, we only use Eq. 12 (replacing \hat{A}_ℓ with the desired method’s attribution matrix in layer ℓ) to calculate the rollout of a given method. However, for methods that do not assume the residual connection, we define a corresponding “FIXED” variation using Eq. 13 that incorporates a fixed value for the context mixing ratio ($r_i = 0.5$).

4 Experiments

In this section, we introduce the datasets and the token attribution analysis methods used in our evaluations, followed by the experimental setup and results.

4.1 Datasets

All analysis methods are evaluated on three different classification tasks. To cover sentiment detection tasks we use SST2 (Socher et al., 2013), MNLI (Williams et al., 2018) for Natural Language Inference and Hatexplain (Mathew et al., 2021) in hate speech detection.

4.2 Analysis Methods

We use two groups explainability approaches in our work: *Weight-based* and *Norm-based*.³ The *Weight-based* approaches employed in our experiments are as follows:

- \mathcal{W} : The raw attention maps averaged across all heads (See \bar{A}_ℓ in §2).
- $\mathcal{W}_{\text{FIXEDRES}}$: Abnar and Zuidema’s assumption; add an identity matrix as a fixed residual to \bar{A}_ℓ (See \hat{A}_ℓ in Eq. 13).
- \mathcal{W}_{RES} : To correct the \mathcal{W} with only the accurate residuals, add the residual based on the context-mixing ratios of \mathcal{N}_{ENC} :

$$\hat{r}_i = \frac{\left\| \sum_{j=1, j \neq i}^n \tilde{\mathbf{x}}_{i \leftarrow j} \right\|}{\left\| \sum_{j=1, j \neq i}^n \tilde{\mathbf{x}}_{i \leftarrow j} \right\| + \|\tilde{\mathbf{x}}_{i \leftarrow i}\|} \quad (14)$$

$$A'_\ell = \mathbf{diag}(\hat{r}_1, \dots, \hat{r}_n) \bar{A}_\ell + \mathbf{diag}(1 - \hat{r}_1, \dots, 1 - \hat{r}_n) I$$

The *Norm-based* analysis methods, namely \mathcal{N} , \mathcal{N}_{RES} and $\mathcal{N}_{\text{RESLN}}$ are discussed in detail in §2. Our proposed norm-based method \mathcal{N}_{ENC} is discussed in §3. For our ablation study, we introduce $\mathcal{N}_{\text{FIXEDRES}}$ which is \mathcal{N} , corrected with a fixed residual similar to $\mathcal{W}_{\text{FIXEDRES}}$.⁴

$$\hat{\mathcal{N}} = \left(\frac{\|z_{i \leftarrow j}\|}{\sum_j \|z_{i \leftarrow j}\|} \right) \in \mathbb{R}^{n \times n} \quad (15)$$

$$\mathcal{N}_{\text{FIXEDRES}} := 0.5 \hat{\mathcal{N}} + 0.5 I$$

We refer to our proposed global method—aggregated \mathcal{N}_{ENC} by the rollout method at the final layer—as *GlobEnc*.

4.3 Gradient-based Methods for Faithfulness Analysis

Gradient-based methods are widely used as alternatives for attention-based counterparts for quantifying the importance of a specific input feature in making the right prediction (Li et al., 2016; Atanasova et al., 2020). In this section we discuss the specific gradient-based methods we use, namely saliency, HTA, and our adjusted HTA.

4.3.1 Saliency

Gradient-based saliency is based on the gradient of the output (y_c) w.r.t. the input embeddings (e_i^0). One of its most accurate variations is the *gradient \times input* method (Kindermans et al., 2016) where the input embeddings is multiplied by the gradients. Thus, the contribution score of input token i is determined by first computing the element-wise product of the input embeddings (e_i^0) and the gradients of the true class output score (y_c) w.r.t. the input embeddings. Then, the L2 norm of the scaled gradients is computed to derive the final score:

$$\text{Saliency}_i = \left\| \frac{\partial y_c}{\partial e_i^0} \odot e_i^0 \right\|_2 \quad (16)$$

³Note that in our experiments, we use all these methods within the rollout aggregation method.

⁴The only difference is that we need to normalize \mathcal{N} before adding an identity matrix.

	Attention Rollout		
	SST2	MNLI	HATEXPLAIN
Weight-based (\mathcal{W})	-0.11 ± 0.26	-0.06 ± 0.22	$+0.12 \pm 0.26$
w/ Fixed Residual ($\mathcal{W}_{\text{FIXEDRES}}$) ⁵	-0.24 ± 0.26	-0.05 ± 0.26	$+0.13 \pm 0.28$
w/ Residual (\mathcal{W}_{RES})	$+0.21 \pm 0.26$	$+0.30 \pm 0.24$	$+0.55 \pm 0.23$
Norm-based (\mathcal{N})	$+0.44 \pm 0.20$	$+0.47 \pm 0.16$	$+0.43 \pm 0.22$
w/ Fixed Residual ($\mathcal{N}_{\text{FIXEDRES}}$)	$+0.48 \pm 0.20$	$+0.55 \pm 0.16$	$+0.48 \pm 0.22$
w/ Residual (\mathcal{N}_{RES})	$+0.73 \pm 0.13$	$+0.75 \pm 0.10$	$+0.66 \pm 0.17$
w/ Residual + Layer Norm 1 ($\mathcal{N}_{\text{RESLN}}$)	-0.21 ± 0.26	-0.06 ± 0.26	$+0.08 \pm 0.28$
w/ GlobEnc : [Residual + Layer Norm 1, 2] (\mathcal{N}_{ENC})	$+0.77 \pm 0.12$	$+0.78 \pm 0.09$	$+0.72 \pm 0.17$

Table 1: Spearman’s rank correlation of attribution based importance (aggregated by rollout) with saliency scores for the validation set for the BERT model fine-tuned on SST-2, MNLI, and HateXplain. In fixed residual cases, the context-mixing ratio is 0.5, and in weight-based w/ residual (\mathcal{N}_{RES}), it is corrected with context-mixing ratio of (\mathcal{N}_{ENC}). The numbers are the average on all the validation set examples \pm the standard deviation.

4.3.2 HTA x Inputs

To determine an upper bound on the information mixing within each layer, we use a modified version of *Hidden Token Attribution* (Brunner et al., 2020, HTA). In the original version, HTA is the sensitivity between any two vectors in the model’s computational graph. However, inspired by the *gradient* \times *input* method (Kindermans et al., 2016), which has shown more faithful results (Atanaseva et al., 2020; Wu and Ong, 2021), we multiply the input vectors by the gradients and then apply a Frobenius norm. We compute the attribution from hidden embedding j ($e_j^{\ell-1}$) to hidden embedding i (e_i^ℓ) in layer ℓ as:

$$c_{i \leftarrow j}^\ell = \left\| \frac{\partial e_i^\ell}{\partial e_j^{\ell-1}} \odot e_j^{\ell-1} \right\|_F \quad (17)$$

Computing HTA-based attribution matrices is an extremely computationally intensive task (especially for long texts) due to the high dimensionality of the hidden embeddings. Hence, we only use this method for 256 examples from the SST-2 task’s validation set. It is worth noting that extracting the HTA-based contribution maps for the aforementioned data took approximately 2 hours, whereas computing the maps for the entire analysis methods stated in §4.2 took only 5 seconds.⁶

⁵As mentioned in §4.2, this analysis method is based on the original experiment by Abnar and Zuidema (2020). Our experiments on SST2 differ from theirs in two aspects: (i) We opted for *gradient* \times *input* saliencies, while they used the sum of gradients (sensitivity) (ii) Instead of BERT, they used a DistillBERT fine-tuned model (Sanh et al., 2019). However, it still yields lower results (Spearman Corr. = 0.13)

⁶Conducted on a 3070 GPU machine.

4.4 Setup

We employ HuggingFace’s transformers library⁷ (Wolf et al., 2020) and the BERT-base-uncased model. For fine-tuning BERT, epochs vary from 3 to 5, and the batch size and learning rate are 32 and 3e-5, respectively.⁸

After rollout aggregation of each analysis method, we obtain an accumulated attribution matrix for every layer (ℓ) of BERT. These matrices indicate the overall contribution of each input token to all token representations in layer ℓ . Since the classifier in a fine-tuned model is attached to the final layer representation of the [CLS] token, we consider the first row (corresponding to [CLS] attributions) of the last layer attribution matrix. This vector represents the contribution of each input token to the model’s final decision. As a measure of faithfulness of the resulting vector with the saliency scores, we report the *Spearman’s rank correlation* between the two vectors.

4.5 Results

Table 1 shows the Spearman correlation of saliency scores with the aggregated attribution scores from [CLS] to input tokens at the final layer. In order to determine the contribution of each component of encoder layer to the overall performance, we report results for multiple attribution analysis methods. Our results demonstrate that incorporating the vector norms, residual connection, and both layer normalizations yields the highest correlation (\mathcal{N}_{ENC}). In what follows next, we discuss the im-

⁷<https://github.com/huggingface/transformers>

⁸Recommended by Devlin et al. (2019).

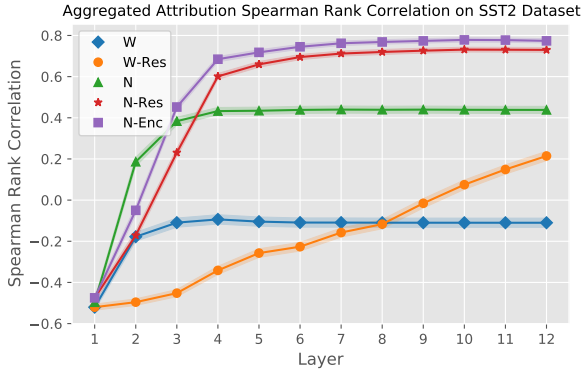


Figure 3: Spearman’s rank correlation of aggregated attribution scores with saliency scores across layers. The 99% confidence intervals are shown as (narrow) shaded areas around each line. \mathcal{N}_{ENC} achieves the highest correlation in almost every layer.

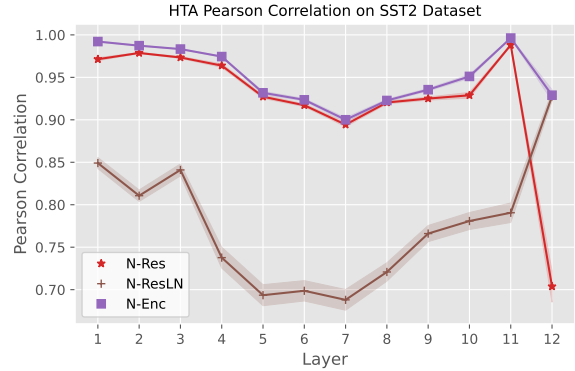


Figure 4: Single layer Pearson correlation of HTA maps with attribution maps. The 99% confidence intervals are shown as shaded areas around each line. $\mathcal{N}_{\text{RESLN}}$ shows considerably less association with HTA.

337 pact of incorporating various parts in the analysis.

338 4.5.1 On the role of vector norms

339 As also suggested by Kobayashi et al. (2020), vec-
 340 tor norms play an important role in determining
 341 attention outputs. This is highlighted by the signif-
 342 icant gap between weight-based and norm-based
 343 settings across all datasets in Table 1.

344 We also show the correlation of the aggregated
 345 attention for all layers in Figure 3. The norm-based
 346 settings (\mathcal{N} and \mathcal{N}_{RES}) attain higher correlation
 347 than the weight-based counterparts (\mathcal{W} and \mathcal{W}_{RES})
 348 almost in all layers, confirming the importance of
 349 incorporating vector norms.

350 4.5.2 On the role of residual connections

351 Kobayashi et al. (2021) showed that in the encoder
 352 layer, the output representations of each token is
 353 mainly determined by its own representation, and
 354 the contextualization from other tokens’ plays a
 355 marginal role. This is in contrary to the simplifying
 356 assumption made by Abnar and Zuidema (2020)
 357 who used a fixed context-mixing ratio of 0.5 (as-
 358 suming that BERT equally preserves and mixes the
 359 representations). This setting is shown as weight-
 360 based with fixed residual ($\mathcal{W}_{\text{FIXEDRES}}$) in Table 1.
 361 We compare this setting against \mathcal{W}_{RES} (see §4.2).
 362 \mathcal{W}_{RES} is similar to $\mathcal{W}_{\text{FIXEDRES}}$ (in that it does not
 363 take into account vector norms) but differs in that
 364 it considers a dynamic mixing ratio (the one from
 365 \mathcal{N}_{ENC}). The huge performance gap between the
 366 two settings in Table 1 clearly highlights the im-
 367 portance of considering accurate context-mixing
 368 ratios. Therefore, it is crucial to consider the resid-

369 ual connection in the attention block for input token
 370 attribution analysis.

371 To further demonstrate the role of residual con-
 372 nections, we utilize the introduced method in §4.2,
 373 where we corrected the norm-based attentions with
 374 fixed residual ($r = 0.5$). The comparison of norm-
 375 based without any residual (\mathcal{N}) and with a fixed
 376 residual ($\mathcal{N}_{\text{FIXEDRES}}$) shows a consistent improve-
 377 ment for the latter across all the datasets. This
 378 provides evidence on that having a fixed uniform
 379 context-mixing ratio is better than neglecting the
 380 residual connection altogether.

381 Finally, when we aggregate the norm-based anal-
 382 ysis with an accurate dynamic context-mixing ratio
 383 (\mathcal{N}_{RES}), we observe the highest correlation up to
 384 this point, without layer normalization.

385 4.5.3 On the role of layer normalization

386 In Table 1 we see a sudden drop in correlations for
 387 $\mathcal{N}_{\text{RESLN}}$. Although this method considers vector
 388 norms and residuals, incorporating LN #1 in the
 389 encoder seems to have deteriorated the accuracy for
 390 token attribution analysis. To determine whether
 391 this deterioration of correlation in aggregated attri-
 392 butions is also present in individual single layers,
 393 we compare the HTA maps as a baseline with the
 394 attribution matrices extracted from different anal-
 395 ysis methods. Figure 4 shows the correlation of
 396 HTA attribution maps with the maps obtained by
 397 \mathcal{N}_{RES} , $\mathcal{N}_{\text{RESLN}}$, and \mathcal{N}_{ENC} methods. The results
 398 indicate that $\mathcal{N}_{\text{RESLN}}$ exhibits a significantly lower
 399 association.

400 The question that arises here is that how incor-
 401 porating an additional component of the encoder
 402 (LN #1) in $\mathcal{N}_{\text{RESLN}}$ degrades the results (compared

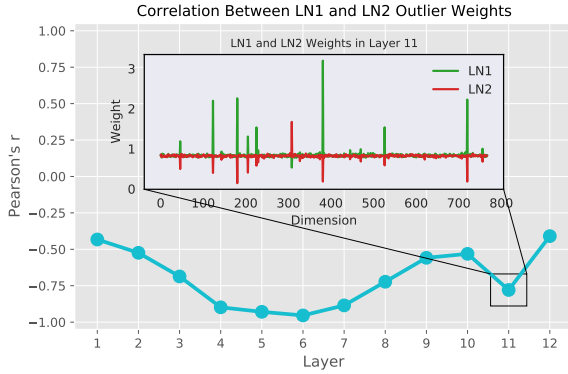


Figure 5: The Pearson correlation between outlier weights of layer normalization #1 and #2 across layers. The exact weights for layer 11 are shown as well.

to \mathcal{N}_{RES}). To answer this question, we investigate the learned weights of layer norm #1 and #2. The outlier weights⁹ in specific dimensions of LNs are shown to be significantly influential on the model’s performance (Kovaleva et al., 2021; Luo et al., 2021). It is interesting to note that based on our observations, the outlier weights of the two layer norms seem to be the opposite of each other. Figure 5 demonstrates the exact weights in layer 11 and also the correlation of the outlier weights across layers. The large negative correlations confirm that the outlier weights work contrary to each other. We speculate that the effect of outliers in the two layer norms is partly cancelled out when both are considered.

As shown in Figure 2, the FFN and the second layer normalization are on top of the attention block. However, $\mathcal{N}_{\text{RESLN}}$ does not incorporate the components outside of the attention block. As described in §3, in our local analysis method \mathcal{N}_{ENC} we incorporate the second layer normalization in the transformer’s encoder (Figure 2), thus considering the whole encoder block (except FFN). Overall, our global method noted as GlobEnc yields the best results among all the methods evaluated in our experiments. In general, Table 1 suggests that incorporating each component of the encoder will increase the correlation; however, the two layer normalizations should be considered together.

4.5.4 On the role of aggregation

We carried out an additional analysis to verify if incorporating vector norms, residual connection and layer normalizations in individual layers is ade-

⁹We identify the dimensions where the weights are at least 3σ from the mean as outliers (Kovaleva et al., 2021).

		L1	L6	L12	MAX
Indiv.	\mathcal{N}	$-.50 \pm .18$	$+.28 \pm .23$	$+.40 \pm .21$	$+.41 \pm .21$
	\mathcal{N}_{RES}	$-.48 \pm .18$	$+.29 \pm .24$	$+.41 \pm .19$	$+.41 \pm .19$
	\mathcal{N}_{ENC}	$-.47 \pm .18$	$+.29 \pm .24$	$+.41 \pm .19$	$+.41 \pm .19$
Rollout	\mathcal{N}	$-.50 \pm .18$	$+.44 \pm .20$	$+.44 \pm .20$	$+.44 \pm .20$
	\mathcal{N}_{RES}	$-.48 \pm .18$	$+.70 \pm .14$	$+.73 \pm .13$	$+.73 \pm .13$
	\mathcal{N}_{ENC}	$-.47 \pm .18$	$+.74 \pm .14$	$+.77 \pm .12$	$+.78 \pm .12$

Table 2: Spearman’s rank correlation of attribution-based scores (individual and aggregated by rollout) with saliency scores for the validation set for the BERT model fine-tuned on SST-2. The results are reported for layers 1, 6, 12, and the maximum of all layers. Rollout aggregation achieves the highest correlations.

quate for achieving high correlations, or if it is also necessary to aggregate them via rollout. Table 2 shows the correlation results in different layers for raw attributions (without aggregation) and for the aggregated attributions using the rollout method. Applying rollout method on attribution maps up to each layer results in higher correlations with the saliency scores than the raw single layer attribution maps, especially in deeper layers. Therefore, attention aggregation is essential for global input token attribution analysis.

An interesting point in Figure 3, which shows the correlation of the aggregated methods throughout the layers, is that the correlation curves flatten out after only a few layers.¹⁰ This indicates that BERT identifies decisive tokens only after the first few layers. The final layers only make minor adjustments to this order. Nevertheless, it is worth noting that the order of attribution does not necessarily imply the model’s final decision and the final result may still change for the better or worse (Zhou et al., 2020).

4.5.5 Qualitative analysis

To qualitatively answer if the aggregated attribution maps provide plausible and meaningful interpretations, we take a closer look at the attribution maps generated by GlobEnc. Figure 1 shows the GlobEnc attribution of the model trained on SST-2. Each layer demonstrates the [CLS] token’s aggregated attribution to input tokens up to the corresponding layer. The example inputs are “people cinema at its finest.” and “big fat waste of time.”, both correctly classified by the model. In both cases, GlobEnc focuses on the relevant words for

¹⁰ \mathcal{W}_{RES} is the only exception with a constant increase; this method is gradually and artificially corrected by \mathcal{N}_{ENC} context mixing ratios.

470 sentiment classification, i.e., “finest” and “waste”.
471 An interesting observation in Figure 1 is that in
472 the first few layers, the [CLS] token mostly at-
473 tends to itself while other tokens have marginal
474 impact. As the representations get more contex-
475 tualized in deeper layers, the attribution correctly
476 shifts to the words which indicate the sentiment
477 of the sentence.¹¹ More examples are shown in
478 Figure A.1. Our qualitative analysis suggests that
479 GlobEnc can be useful for a reasonable interpreta-
480 tion of attention mechanism in BERT and possibly
481 any other transformer-based model.

482 5 Related Work

483 While numerous studies have used attention
484 weights to analyze and interpret the self-attention
485 mechanism (Clark et al., 2019; Kovaleva et al.,
486 2019; Reif et al., 2019; Htut et al., 2019), the use
487 of mere attention weights to explain a model’s in-
488 ner workings has been an active topic of debate
489 (Serrano and Smith, 2019; Jain and Wallace, 2019;
490 Wiegrefe and Pinter, 2019). Several solutions have
491 been proposed to address this issue, usually through
492 converting raw attention weights to scores that pro-
493 vide better explanations. Brunner et al. (2020) used
494 the transformation function $f^h(x_j)$ to introduce
495 *effective attentions*—the orthogonal component of
496 the attention matrix in $f^h(x_j)$ null space—to ex-
497 plain the inner workings of each layer. However,
498 this technique ignores other components in the en-
499 coder and is computationally expensive due to the
500 SVD required to compute the effective attentions.
501 Kobayashi et al. (2020) incorporated the modified
502 vector and introduced a vector norms-based analy-
503 sis. This was later extended by integrating residual
504 connections and layer normalization components to
505 enhance the accuracy of explanations (Kobayashi
506 et al., 2021). But, as discussed in §4.5, relying
507 solely on LN #1 does not produce accurate results.

508 While these methods can be employed for single-
509 layer (local) analysis, multi-layer attributions are
510 not necessarily correlated with single-layer attribu-
511 tions due to the significant degree of information
512 combination through multi-layer language mod-
513 els (Pascual et al., 2021; Brunner et al., 2020).
514 Various saliency methods exist for explaining the
515 model’s decision based on the input (Li et al., 2016;
516 Bastings and Filippova, 2020; Atanasova et al.,
517 2020; Wu and Ong, 2021; Mohebbi et al., 2021).

¹¹Complete attention maps in Figure A.2 show that, simi-
larly to [CLS], other tokens also focus on sentiment tokens.

518 However, these approaches are not primarily de-
519 signed for computing inter-token attributions. To
520 fill this gap, Brunner et al. (2020) proposed HTA,
521 which is based on the gradient of each hidden em-
522 bedding in relation to the input embeddings. In
523 §4.3.2, we extend HTA to incorporate the impact
524 of the input vectors. However, HTA is extremely
525 computationally intensive. Attention rollout (see
526 §3) and attention flow—which involve solving a
527 max-flow problem on the attention graph—are two
528 aggregation approaches introduced by Abnar and
529 Zuidema (2020), in which raw attention weights
530 (with equally weighted residual weights) are ag-
531 gregated within multiple layers. We showed that
532 attention rollout does not perform well on a BERT
533 model fine-tuned on downstream tasks and that this
534 problem can be resolved by utilizing attribution
535 norms.

536 6 Conclusions

537 In this work, we proposed a novel method for single
538 layer token attribution analysis which incorporates
539 the whole encoder layer, i.e., the attention block
540 and the output layer normalization. When aggre-
541 gated across layers using the rollout method, our
542 technique achieves quantitatively and qualitatively
543 plausible results. Our evaluation of different analy-
544 sis methods provided evidence on roles played by
545 individual components of the encoder layer, i.e.,
546 the vector norms, the residual connections, and the
547 layer normalizations. Furthermore, our in-depth
548 analysis suggested that the two layer normaliza-
549 tions in the encoder layer counteract each other;
550 hence, it is important to couple them for an accu-
551 rate analysis.

552 Additionally, using a newly proposed and im-
553 proved version of Hidden Token Attribution, we
554 demonstrated that encoder-based attribution analy-
555 sis is more accurate when compared to other partial
556 solutions in a single layer (local-level). This is con-
557 sistent with our global observations. Quantifying
558 global input token attribution based on our work
559 can provide a meaningful explanation of the whole
560 model’s behavior. In future work, one can apply
561 our global analysis method on various datasets and
562 models, to provide valuable insights into model
563 decisions and interpretability.

564
565
566
567
568
569

570
571

572
573
574
575
576
577
578

579
580
581
582
583
584
585

586
587
588
589
590

591
592
593
594
595
596
597

598
599
600
601
602
603
604
605
606

607
608
609
610

611
612
613
614
615
616

617
618
619
620

References

Samira Abnar and Willem Zuidema. 2020. [Quantifying attention flow in transformers](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4190–4197, Online. Association for Computational Linguistics.

Jay Alammar. 2018. [The illustrated transformer \[blog post\]](#).

Pepa Atanasova, Jakob Grue Simonsen, Christina Lioma, and Isabelle Augenstein. 2020. [A diagnostic study of explainability techniques for text classification](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3256–3274, Online. Association for Computational Linguistics.

Jasmijn Bastings and Katja Filippova. 2020. [The elephant in the interpretability room: Why use attention as explanation when we have saliency methods?](#) In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 149–155, Online. Association for Computational Linguistics.

Gino Brunner, Yang Liu, Damian Pascual, Oliver Richter, Massimiliano Ciaramita, and Roger Wattenhofer. 2020. [On identifiability in transformers](#). In *International Conference on Learning Representations*.

Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. [What does BERT look at? an analysis of BERT’s attention](#). In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Phu Mon Htut, Jason Phang, Shikha Bordia, and Samuel R. Bowman. 2019. [Do attention heads in BERT track syntactic dependencies?](#) *CoRR*, abs/1911.12246.

Sarthak Jain and Byron C. Wallace. 2019. [Attention is not Explanation](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3543–3556, Minneapolis, Minnesota.

Pieter-Jan Kindermans, Kristof Schütt, Klaus-Robert Müller, and Sven Dähne. 2016. [Investigating the influence of noise and distractors on the interpretation of neural networks](#).

Goro Kobayashi, Tatsuki Kuribayashi, Sho Yokoi, and Kentaro Inui. 2020. [Attention is not only a weight: Analyzing transformers with vector norms](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7057–7075, Online. Association for Computational Linguistics.

Goro Kobayashi, Tatsuki Kuribayashi, Sho Yokoi, and Kentaro Inui. 2021. [Incorporating Residual and Normalization Layers into Analysis of Masked Language Models](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4547–4568, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Olga Kovaleva, Saurabh Kulshreshtha, Anna Rogers, and Anna Rumshisky. 2021. [BERT busters: Outlier dimensions that disrupt transformers](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3392–3405, Online. Association for Computational Linguistics.

Olga Kovaleva, Alexey Romanov, Anna Rogers, and Anna Rumshisky. 2019. [Revealing the dark secrets of BERT](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4365–4374, Hong Kong, China.

Jiwei Li, Xinlei Chen, Eduard Hovy, and Dan Jurafsky. 2016. [Visualizing and understanding neural models in NLP](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 681–691, San Diego, California. Association for Computational Linguistics.

Ziyang Luo, Artur Kulmizev, and Xiaoxi Mao. 2021. [Positional artefacts propagate through masked language model embeddings](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5312–5327, Online. Association for Computational Linguistics.

Binny Mathew, Punyajoy Saha, Seid Muhie Yimam, Chris Biemann, Pawan Goyal, and Animesh Mukherjee. 2021. [Hatexplain: A benchmark dataset for explainable hate speech detection](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(17):14867–14875.

Hosein Mohebbi, Ali Modarressi, and Mohammad Taher Pilehvar. 2021. [Exploring the role of BERT token representations to explain sentence probing results](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 792–806, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Damian Pascual, Gino Brunner, and Roger Wattenhofer. 2021. [Telling BERT’s full story: from local attention](#)

621
622
623
624
625
626
627

628
629
630
631
632
633
634
635

636
637
638
639
640
641

642
643
644
645
646
647
648

649
650
651
652
653
654
655

656
657
658
659
660
661
662
663

664
665
666
667
668
669

670
671
672
673
674
675
676

677
678

679	to global aggregation. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 105–124, Online. Association for Computational Linguistics.	736
680		737
681		
682		
683		
684	Emily Reif, Ann Yuan, Martin Wattenberg, Fernanda B Viegas, Andy Coenen, Adam Pearce, and Been Kim. 2019. Visualizing and measuring the geometry of bert. In <i>Advances in Neural Information Processing Systems</i> , pages 8594–8603.	
685		
686		
687		
688		
689	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. In <i>NeurIPS EMC² Workshop</i> .	
690		
691		
692		
693	Sofia Serrano and Noah A. Smith. 2019. Is attention interpretable? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 2931–2951, Florence, Italy.	
694		
695		
696		
697	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In <i>Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing</i> , pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.	
698		
699		
700		
701		
702		
703		
704		
705	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.	
706		
707		
708		
709		
710	Sarah Wiegrefe and Yuval Pinter. 2019. Attention is not not explanation. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 11–20, Hong Kong, China. Association for Computational Linguistics.	
711		
712		
713		
714		
715		
716		
717	Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.	
718		
719		
720		
721		
722		
723		
724		
725		
726	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	738
727		739
728		740
729		741
730		
731		
732		
733		
734		
735		
	Zhengxuan Wu and Desmond C. Ong. 2021. On explaining your explanations of BERT: an empirical study with sequence classification. <i>CoRR</i> , abs/2101.00196.	
	Wangchunshu Zhou, Canwen Xu, Tao Ge, Julian McAuley, Ke Xu, and Furu Wei. 2020. Bert loses patience: Fast and robust inference with early exit. In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 18330–18341. Curran Associates, Inc.	742
		743
		744
		745
		746
		747
	A Appendix	748
	A.1 LN Formulation	749
	$m(\mathbf{a}) := \frac{1}{d} \sum_k \mathbf{a}^{(k)}$,	750
		751
	$s(\mathbf{a}) := \sqrt{\frac{1}{d} \sum_k (m(\mathbf{a}) - \mathbf{a}^{(k)} + \epsilon)^2}$	752
	where ϵ is a small constant	753
	A.2 More examples	754
	Aggregated attributions by different methods throughout layers is shown in Figure A.1. Our proposed method shows more plausible results.	755
		756
		757
	Aggregated attribution map for layer 12 is shown in Figure A.2. In this figure, the effect of each token can be seen on all other tokens and not just the [CLS] token.	758
		759
		760
		761

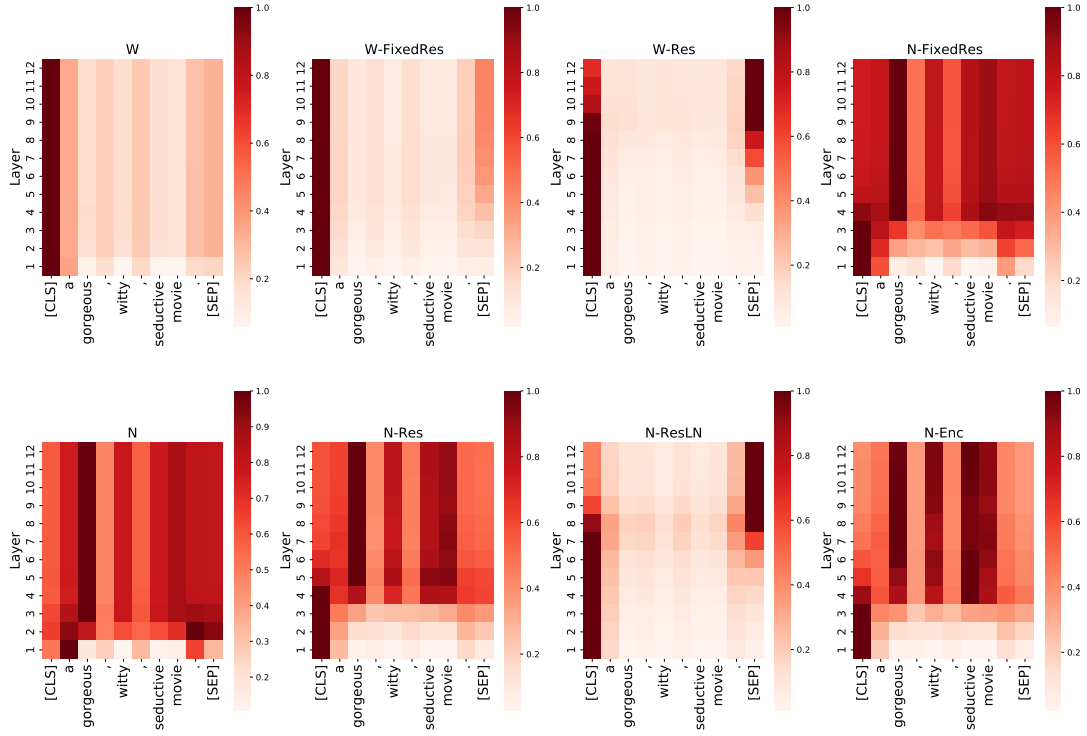


Figure A.1: Spearman correlation for aggregated attributions via rollout with different methods across layers. The model is fine-tuned on SST2 dataset and the attention of the CLS token is shown in each layer.

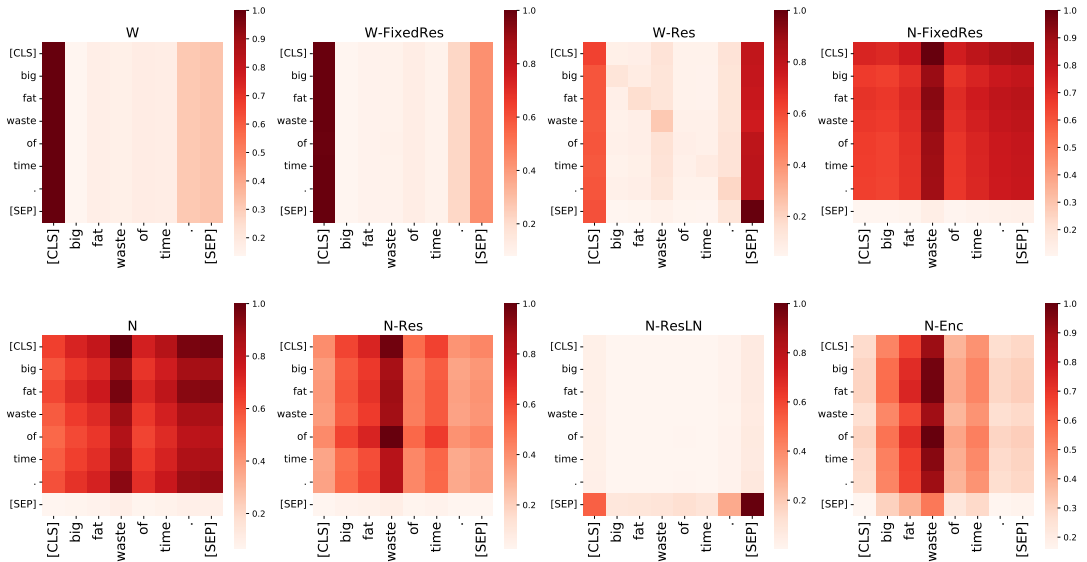


Figure A.2: Spearman correlation for aggregated attributions via rollout with different methods in layer 12. The model is fine-tuned on SST2 dataset. Each row indicates how much other tokens impact the token written on the row.