

# Graph-tree Fusion Model with Bi-Directional Information Propagation for Long Document Classification

Anonymous ACL submission

## Abstract

Long document classification presents challenges in capturing both local and global dependencies due to their extensive content and complex structure. Existing methods often struggle with token limits and fail to adequately model hierarchical relationships within documents. To address these constraints, we propose a novel model leveraging a graph-tree structure. Our approach integrates syntax trees for sentence encodings and document graphs for document encodings, which capture fine-grained syntactic relationships and broader document contexts. We use Tree Transformers to generate sentence encodings, while a graph attention network models inter- and intra-sentence dependencies. During training, we implement bi-directional information propagation from word-to-sentence-to-document and vice versa, which enriches the contextual representation. Our proposed method enables a comprehensive understanding of content at all hierarchical levels and effectively handles arbitrarily long contexts without token limit constraints. Experimental results demonstrate the effectiveness of our approach in all types of long document classification tasks.

## 1 Introduction

Long document understanding has garnered increasing attention in the field of natural language processing (NLP) due to its wide range of applications across various domains, including legal document analysis, scientific literature categorization, and clinical text mining. Accurate understanding of long documents is essential for tasks such as information retrieval, content summarization, and decision-making support systems. Modern deep learning models for semantic analysis achieve impressive results by training on large datasets, which enables them to generate highly accurate predictions on unseen content (Al-Qurishi, 2022). However, their ability to capture relationships between

words and sentences relies on increasingly complex statistical operations as the text sequence lengths (Tay et al., 2020). Consequently, many existing methods become impractical for real-world applications, which makes processing long documents a challenging task.

One of the primary challenges of long document classification is managing the large volume of information. Unlike short texts, long documents contain extensive content that often spans multiple topics, making it difficult to capture the overall context of the document effectively. Transformer-based models (Vaswani et al., 2017), such as BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), GPT-3 (Brown et al., 2020) and LLaMa-2 (Touvron et al., 2023), had gained popularity for NLP tasks due to their ability to capture (relatively) long-range dependencies and contextual relationships. However, in long document classification, transformer models face scalability issues due to their quadratic time complexity with respect to the input length. Processing long documents with transformers can be computationally expensive and memory-intensive, often requiring substantial hardware resources. Current methods for handling lengthy documents include truncating texts to a predefined length or modifying the attention mechanism. Truncating to the first 512 tokens is straightforward but may cause significant information loss. Sparse attention models like Longformer (Beltagy et al., 2020) and Big Bird (Zaheer et al., 2020) reduce computational load by focusing on a subset of tokens, but they do not fully capture comprehensive context and dependencies in long texts.

Another significant challenge in long document classification is capturing contextual dependencies and the hierarchical structure. Long documents have dependencies at word, sentence, and document levels, which are crucial for understanding content and leveraging both local and global contexts. This structure is important for understand-

ing the overall context and meaning. Treating text as a flat sequence of tokens can cause models to miss these important hierarchical relationships. Current approaches to address these challenges involve hierarchical models. For instance, Hierarchical Attention Network (HAN) (Yang et al., 2016) and Hierarchical Attention Transformers (HAT) (Chalkidis et al., 2022), aim to capture both word-level and sentence-level representations before aggregating them into document-level embeddings. However, these models often fail to capture the intricate relationships between different parts of the document, such as the interplay between words, sentences, and overall document structure (Dai et al., 2019). Additionally, hierarchical models may struggle with long-range dependencies (Dong et al., 2023), which misses the relationships between distant sections essential for understanding the overall context.

To address the aforementioned challenges and constraints, we propose a novel model that leverages a graph-tree structure for arbitrarily long document classification. Our approach fuses syntax trees for sentence encodings with document graphs for document encodings, which provides a comprehensive representation that captures both local and global dependencies. Syntax trees (both dependency and constituency) represent the grammatical structure of sentences, which enhances sentence-level understanding. We use Tree Transformers (Ahmed et al., 2019a) to generate sentence encodings from these syntax trees. The document graph preserves hierarchical relationships within the document, which ensures that both local and global contexts are considered during the classification process. We apply the Graph Attention Network (GAT) (Veličković et al., 2018) on the constructed document graph to model the dependencies between the sentences. This graph structure effectively captures both inter- and intra-sentence dependencies. Meanwhile, during training, we implement a bidirectional information propagation approach where information flows both from word-to-sentence-to-document and from document-to-sentence-to-word. This bidirectional flow enriches the contextual representation of the document, which allows for a more comprehensive understanding of the content at all hierarchical levels. By incorporating syntax trees and document graphs, we also encode the semantic units differently according to their characteristics. This allows our model to handle arbitrarily long contexts without

being constrained by token limits.

To summarize, our main contributions are:

1. We introduce a novel graph-tree structure that combines syntax trees and document graphs to capture both local and global dependencies within arbitrarily long documents.
2. We introduce a bidirectional information propagation approach where the information flows both from word-to-sentence-to-document and from document-to-sentence-to-word, which enriches the contextual representation of the document.
3. We show empirically that our model achieves improvements across a variety of classification tasks, including binary, multi-class, and multi-label classification.

## 2 Preliminaries

### 2.1 Tree Transformer

Tree Transformer (Ahmed et al., 2019a) is designed to more effectively preserve syntactic and semantic information. Given a dependency or constituency tree structure of a sentence, a dependency tree has a word at every node, represented by  $\mathbf{X}_d$  while, in a constituency tree, only the leaf nodes contain words, represented by  $\mathbf{X}_c$ :

$$X_d = \begin{bmatrix} \mathbf{p}_v \\ \mathbf{c}_1^v \\ \mathbf{c}_2^v \\ \vdots \\ \mathbf{c}_n^v \end{bmatrix}, \quad X_c = \begin{bmatrix} \mathbf{c}_1^v \\ \mathbf{c}_2^v \\ \vdots \\ \mathbf{c}_n^v \end{bmatrix}, \quad (1)$$

where  $\mathbf{p}_v$  is the initial parent representation and  $\mathbf{c}_i^v$  is the initial child representation of node  $i$ .

The parent node embedding  $\mathbf{P}$  is computed using multi-branch attention built upon multi-head attention in the vanilla Transformer (Vaswani et al., 2017). The branch attention  $\mathbf{B}_i$  for branch  $i$  is computed as :

$$\mathbf{B}_i = \text{Attention}(\mathbf{Q}_i \mathbf{W}_i^Q, \mathbf{K}_i \mathbf{W}_i^K, \mathbf{V}_i \mathbf{W}_i^V), \quad (2)$$

where  $\mathbf{W}_i^Q$ ,  $\mathbf{W}_i^K$ , and  $\mathbf{W}_i^V$  are the learnable weight matrices. Each  $\mathbf{B}_i$  is then normalized and scaled using a layer normalization block (Ba et al., 2016):

$$\bar{\mathbf{B}}_i = \text{LayerNorm}(\mathbf{B}_i \mathbf{W}_i^b + \mathbf{B}_i) \times \kappa_i, \quad (3)$$

where  $\mathbf{W}_i^b$  and  $\kappa_i$  are the learnable parameters. Then, a position-wise convolutional neural network

(PCNN) is applied to each  $\bar{\mathbf{B}}_i$ , and the branch attention is aggregated:

$$\text{BranchAttn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \sum_{i=1}^n \alpha_i \text{PCNN}(\bar{\mathbf{B}}_i), \quad (4)$$

where  $\alpha_i$  is learnable. The final parent representation, or sentence embedding, is obtained by:

$$\mathbf{P}' = \text{EwS}(\tanh((\mathbf{x}' + \mathbf{x})\mathbf{W} + b)), \quad (5)$$

where EwS is element-wise summation, and  $\mathbf{x}$  and  $\mathbf{x}'$  depicts the input and output of the attention module. Additional details on the full operations of the Tree Transformer are in the original paper (Ahmed et al., 2019a).

## 2.2 Graph Attention Network

Graph Attention Network (GAT) (Veličković et al., 2018) is designed to model information flow between nodes, which enhances node representations by employing attention over features from neighbouring nodes. Given a heterogeneous graph  $G = (V, E)$  with  $N$  nodes, the input node features  $h = \{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_N\}$ , and GAT layer with multi-head attention is designed as follows:

$$e_{ij} = \text{LeakyReLU}(\mathbf{a}^T [\mathbf{W}\mathbf{h}_i \parallel \mathbf{W}\mathbf{h}_j]),$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k \in \mathcal{N}_i} \exp(e_{ik})}, \quad (6)$$

$$\mathbf{h}'_i = \parallel_{k=1}^K \sigma \left( \sum_{j \in \mathcal{N}_i} \alpha_{ij}^k \mathbf{W}^k \mathbf{h}_j \right),$$

where  $\parallel$  denotes concatenation,  $\mathbf{W}$  is a weight matrix,  $\mathbf{a}$  is a weight vector,  $\mathcal{N}_i$  is the neighbourhood of node  $i$ , LeakyReLU and  $\sigma$  are the activation functions,  $K$  is the number of attention heads,  $\alpha_{ij}^k$  and  $\mathbf{W}^k$  are the attention coefficients and weight matrix for the  $k^{\text{th}}$  head, respectively.

## 3 Method

We propose a novel Graph-Tree Fusion Model, as shown in Figure 1 that leverages a graph-tree structure for long document classification. Our model uses multi-granularity document representations through Tree Transformers and graph attention networks. It fuses syntax trees for sentence encodings with document graphs for document encodings. Additionally, during training, we implement a bidirectional information propagation approach, allowing information to flow both from word-to-sentence-to-document and from document-to-sentence-to-word.

## 3.1 Sentence Encoder via Syntax Tree

Sentences are foundational units of a document. Preserving sentence semantics enhances the accuracy and informativeness of document embeddings.

We first split a document  $\mathcal{D}$  into sentences as  $\mathcal{S} = \{s_1, s_2, \dots, s_N\}$ . Given a sentence  $s$  with a sequence of input tokens, we use RoBERTa (Liu et al., 2019) to encode the tokens and output the corresponding vector for each token from the last hidden layer, denoted as  $\mathbf{E}(s) = \{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\}$ , where  $\mathbf{e}_i$  is the embedding for token  $i$ . We then parse the sentences, obtaining initialized sentence embeddings  $\mathbf{X}_d^{s_i}$  for the dependency tree and  $\mathbf{X}_c^{s_i}$  for the constituency tree (see equation 1). These two tree structures are then processed by the Tree Transformers mentioned in Section 2.1, yielding enhanced sentence representations for each sentence:  $\mathbf{h}_d^{s_i}$  for the dependency tree and  $\mathbf{h}_c^{s_i}$  for the constituency tree:

$$\begin{aligned} \mathbf{h}_d^{s_i} &= \text{TreeTransformer}(\mathbf{X}_d^{s_i}), \\ \mathbf{h}_c^{s_i} &= \text{TreeTransformer}(\mathbf{X}_c^{s_i}). \end{aligned} \quad (7)$$

The final embedding for sentence  $i$  is defined as the mean-pooling of the two tree representations:

$$\mathbf{h}^{s_i} = \frac{1}{2}(\mathbf{h}_d^{s_i} + \mathbf{h}_c^{s_i}). \quad (8)$$

## 3.2 Document Encoder via Document Graph Sentence Selector Using Label-wise Attention.

To identify the most important sentences, we calculate the similarity score between each sentence  $s_i$  in  $\mathcal{D}$  and the labels using label-wise attention. We obtain the label embeddings for each label  $\mathcal{Y}_i$  by averaging the word embeddings (from RoBERTa) of each word in their descriptors:

$$\mathbf{h}^{\mathcal{Y}_i} = \frac{1}{m} \sum_{j \in m} w_j, i = 1, 2, \dots, L, \quad (9)$$

where  $m$  is the number of words in the descriptor, and  $L$  is the number of labels. The similarity score is:

$$\alpha_{s_i, \mathcal{Y}_i} = \text{Softmax}(\mathbf{h}^{s_i} \cdot \mathbf{h}^{\mathcal{Y}_i}), \quad (10)$$

where  $\alpha_{s_i, \mathcal{Y}_i}$  is a probability that contains the scores associated with a sentence  $i$  in  $\mathcal{D}$  to a specific label  $\mathcal{Y}_i$ . We then apply a threshold  $\tau$  to select sentences with high probabilities:

$$\mathcal{S}' = \{s_i | \alpha_{s_i, \mathcal{Y}_i} \geq \tau\}. \quad (11)$$

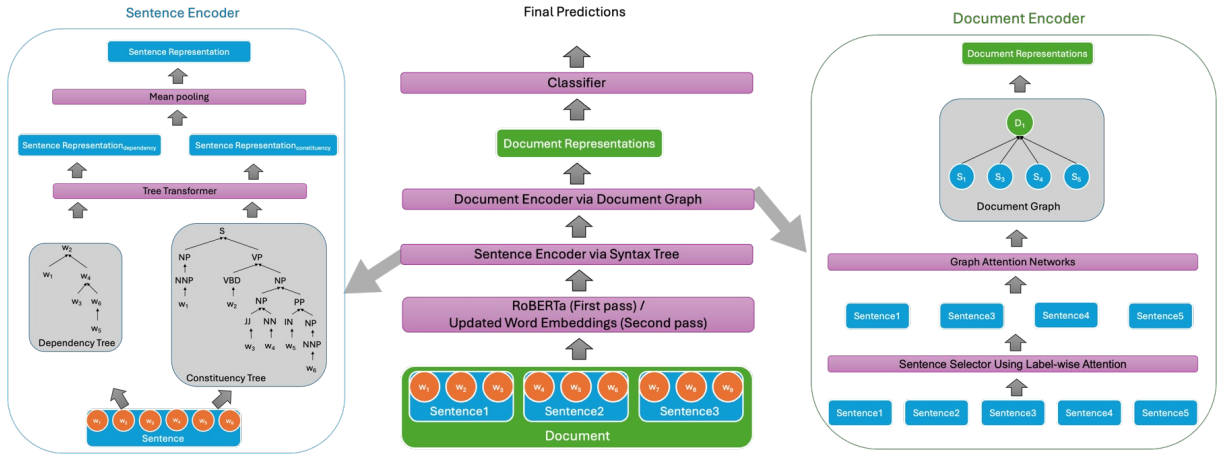


Figure 1: Overview of our proposed model. Our proposed model integrates syntax trees for sentence encodings (showing on the left) and document graphs for document encodings (showing on the right). We employ Tree Transformers to generate sentence encodings from the syntax trees and use graph attention networks to generate document encodings from the document graph. Our model workflow includes two passes. In the first pass, initial word embeddings are obtained from RoBERTa. In the second pass, these word embeddings are updated through bidirectional information propagation.

**Document Encoding Using Graph Attention Network (GAT).** To obtain the document representation, we construct a heterogeneous document graph  $G = (V, E)$  that captures the relations between documents and their sentences. The graph  $G$  contains sentence nodes and document nodes, and one type of edge: document-sentence edges. Specially, for each document  $\mathcal{D}$ , we create a document node  $v_{\mathcal{D}}$  and a set of selected sentence nodes  $V_{S'} = \{v_{s_1}, v_{s_2}, \dots, v_{s_N}\}$ . Directed edges  $E$  are established from each sentence node  $v_{s_i}$  to its corresponding document node  $v_{\mathcal{D}}$ .

We apply GAT (described in 2.2) to model the inter-sentence and document relations within a document, which enhances the document representations derived from sentence representations.

To obtain the document representation for  $\mathcal{D}$ , the sentence node embeddings  $\mathbf{h}^{s_i}$  are initially generated using the sentence encoder. The document node feature is then initialized by taking the mean pooling of the features of its sentence nodes:

$$\mathbf{h}^{\mathcal{D}} = \frac{1}{|V_{S'}|} \sum_{s_i \in V_{S'}} \mathbf{h}^{s_i}. \quad (12)$$

We use the GAT layer with multi-head attention to compute the new document node features  $\mathbf{h}'^{\mathcal{D}}$  as follows:

$$\mathbf{h}'^{\mathcal{D}} = \text{GAT}(\mathbf{h}^{s_i}). \quad (13)$$

After each GAT layer, we introduce a position-wise feed-forward (FFN) layer, consisting of two linear transformations similar to the vanilla Transformer

architecture (Vaswani et al., 2017), to obtain the final document representation:

$$\tilde{\mathbf{h}}^{\mathcal{D}} = \text{FFN}(\mathbf{h}'^{\mathcal{D}}). \quad (14)$$

### 3.3 Bidirectional Information Propagation

Inspired by Wang et al. (2020), we implement a bidirectional information propagation approach, as shown in Figure 2. This approach allows information to flow from word-to-sentence-to-document and from document-to-sentence-to-word.

After obtaining the document representation from the document encoder in Section 2.2, we update the sentence nodes using the updated document nodes and then update the word nodes using the updated sentence encodings. We further iteratively update the document nodes and sentence nodes. The information flow is bidirectional: in each iteration, it first moves from word-to-sentence-to-document (see Sections 3.1,3.2), and then from document-to-sentence-to-word. For the  $t^{\text{th}}$  iteration, the document-to-sentence-to-word process can be represented as:

$$\begin{aligned} \mathbf{U}_{\mathcal{D} \rightarrow \mathcal{S}}^{t+1} &= \text{GAT}(\mathbf{H}_{\mathcal{D}}^t), \\ \mathbf{H}_{\mathcal{S}}^{t+1} &= \text{FFN}(\mathbf{U}_{\mathcal{D} \rightarrow \mathcal{S}}^{t+1}), \\ \mathbf{U}_{\mathcal{S} \rightarrow w}^{t+1} &= \text{GAT}(\mathbf{H}_{\mathcal{S}}^{t+1}), \\ \mathbf{H}_w^{t+1} &= \text{FFN}(\mathbf{U}_{\mathcal{S} \rightarrow w}^{t+1}) \end{aligned} \quad (15)$$

where  $\mathbf{H}_{\mathcal{D}}$  is initialized by  $\tilde{\mathbf{h}}^{\mathcal{D}}$  in Equation 3.2 in the first iteration.

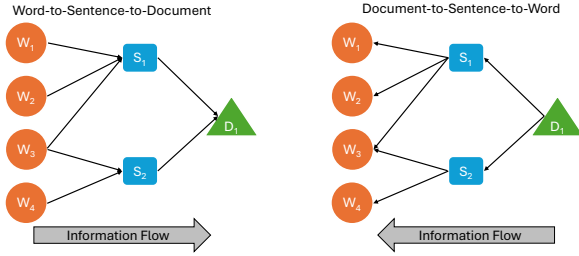


Figure 2: The detailed bidirectional information propagation. Orange, blue, and green nodes represent word, sentence, and document nodes, respectively. The arrows on the edges indicate the current direction of information flow. First, on the left, words are used to aggregate sentence-level information, and the resulting sentence representations are then used to aggregate document-level information. Next, on the right, sentences are updated with the new document representations, and words are updated with the new sentence representations.

As illustrated in Figure 2, word nodes can aggregate document-level information from sentences. For example, a word node with a high degree indicates frequent occurrences in multiple sentences, suggesting it is a keyword of the document. Sentence nodes with a higher concentration of important words are more likely to contain significant information, making them suitable for forming key sections. Bidirectional information propagation enables a comprehensive exchange of information across different hierarchical levels.

### 3.4 Model Workflow

The architectural structures of the proposed models are portrayed in Figure 1. The proposed model requires two forward passes separated by document-to-sentence-to-word update step. During the initial forward pass, RoBERTa word embeddings are used as the initial input and undergo simultaneous processing by both the Dependency Tree Transformers (DTT) and Constituency Tree Transformers (CTT), and finally mean-pooling (see formula 8) in the sentence encoder. The GAT layer in the document encoder computes the document representation utilizing the sentence representations from the sentence encoder and the first forward pass ends here.

The completion of the first forward pass initiates the activation of the document-to-sentence-to-word update step. The initial step employs the document-to-sentence update step which updates the the  $\mathbf{h}_d^{s_i}$  ( $\mathbf{h}_d^{s_i}$ ) and  $\mathbf{h}_c^{s_i}$  ( $\mathbf{h}_c^{s_i}$ ). Then, the sentence-to-word refinement step is utilized twice: once to update the word embeddings based on the updated  $\mathbf{h}_d^{s_i}$ ,

Dataset	Type	# of Classes	# of Instance	Average # of Words per Document
Hyperpartisan	Binary	2	754,000	745
AMZ	Multi-class	5	4850	12,356
20News	Multi-class	20	20,000	369
BOOK		227	16,559	575
ECtHR	Multi-label	33	11,000	5530
Essays		5	1255	660

Table 1: Details of the Long Document Classification Datasets.

and another time based on the updated  $\mathbf{h}_c^{s_i}$ .

After the document-to-sentence-to-word update step is employed, the second forward pass is initiated. It works over the pruned syntax trees and graph representation of sentences and document nodes given by the first forward pass. This second pass is almost identical to the first pass with a small twist. Here, the sentence encoder works with two different word embeddings. The CTT intakes the word embeddings updated by the  $\mathbf{h}_c^{s_i}$ , and the DTT is fed with word embeddings updated by the  $\mathbf{h}_d^{s_i}$  as inputs. The following steps work in the similar manner to the first forward pass. This second forward pass generates refined sentence representations ( $\mathbf{h}_c^{s_i}$ , and  $\mathbf{h}_d^{s_i}$ ) and subsequently the refined document representation ( $\mathbf{H}'_D$ ).

## 4 Experiment

### 4.1 Setup

**Datasets.** We evaluate our proposed model on six common long document classification datasets: CMU BOOK Summary (Bamman and Smith, 2013), ECtHR (Chalkidis et al., 2021), Hyperpartisan (Kiesel et al., 2019), 20News (Lang, 1995), Amazon product reviews (AMZ) (He and McAuley, 2016), and Essays (Pennebaker and King, 1999). Following Lu et al. (2023), we randomly sample product reviews longer than 2048 words from the Book category for the AMZ dataset. The statistics of the datasets are summarized in Table 4.1.

**Implementation Details.** The model employs an initial learning rate of 0.1, which is subsequently reduced by 80% in each iteration if the validation accuracy declines compared to the previous iteration. The batch size is 10. For the tree-transformers, the same hyper-parameter settings are used as in Ahmed et al. (2019b). The statement encoding unit utilizes a GAT with six attention heads. The model’s parameters are trained using the “Adam” optimizer (Lydia and Francis, 2019). The performance evaluation of our models has been conducted using 10-fold cross-validation. To facili-

tate this cross-validation process, we have utilized the StratifiedKFold function from the scikit-learn package. All experiments have been conducted in an Ubuntu 22.04 LTE environment, leveraging a 48GB NVIDIA RTX A6000 GPU. For parsing the sentences and generating the tree representations, we have used the Stanford Core-NLP parser (Manning et al., 2014).

## 4.2 Baseline Models

Following Lu et al. (2023), we compare our methods with Transformer-based model with pre-training and State-Space Model (SSM) system.

**BERT with Pre-training** This simplest approach involves fine-tuning BERT (Devlin et al., 2019) after truncating long documents to the first 512 tokens. A fully connected layer is then applied to the [CLS] token for classification.

**ToBERT** Transformer over BERT (ToBERT) is a hierarchical approach designed to process documents of any lengths (Pappagari et al., 2019). It divides long documents into chunks of 200 tokens and applies a Transformer layer to the BERT-based representations of these chunks.

**Longformer** is designed to handle longer input sequences with efficient self-attention that scales linearly with the sequence length, allowing it to process up to 4,096 tokens (Beltagy et al., 2020).

**BERT+TextRank** To address BERT’s 512-token limitation, Park et al. (2022) augment the first 512 tokens with a second set of 512 tokens selected using TextRank (Mihalcea and Tarau, 2004), an efficient unsupervised sentence ranking algorithm.

**BERT+Random** As an alternative method to BERT+TextRank, Park et al. (2022) augment the first 512 tokens by selecting random sentences up to an additional 512 tokens.

**Hungry Hungry Hippo with Max Pooling (H3-pooler)** H3 (Fu et al., 2023) is an SSM-based method designed for simultaneous multi-object tracking, maintaining and updating object states based on observed data. Lu et al. (2023) enhanced this model by inserting a max pooling layer between each SSM block, creating the H3-Pooler.

## 5 Results and Discussions

### 5.1 Performance Comparison

We compare our proposed model against previous baseline models on various evaluation metrics, as shown in Table 2. We report accuracy for binary

and multi-class tasks (Hyperpartisan, 20News, and AMZ) and macro-F1 scores for the multi-label classification problems (BOOK, ECtHR, and Essays). Each row in the table presents the performance of a specific method on each dataset, with the best score for each dataset highlighted.

The results demonstrate that our model consistently achieves superior performance across all datasets. For binary classification (Hyperpartisan), our model outperforms existing methods by approximately 1.3%. In multi-class classification tasks (20News and AMZ), our model shows a marked improvement, exceeding baseline accuracies by about 2% to 4%. In multi-label classification tasks (BOOK, ECtHR, and Essays), our model demonstrates a significant enhancement, with macro-F1 scores improving by approximately 3% to 10% over the best baseline models. These results underscore the robustness and effectiveness of our proposed model in handling the complexities of long document classification across different types of classification tasks.

### 5.2 Ablation Studies

We are interested in studying the effectiveness and robustness of our model by analyzing various components, such as the tree structure, graph structure, and bidirectional information propagation. To understand the impacts of these factors, we conduct controlled experiments with three different settings: (a) removing the tree structure, which is further divided into three sub-experiments: removing the CTT structure, removing the DTT structure, and replacing the entire tree structure with the [CLS] token; (b) replacing the GAT with a max-pooling layer; and (c) removing the bidirectional information propagation. This allows us to evaluate the influence of each module individually, without interference from the others. The results are summarized in Table 3.

**Effectiveness of the Tree Structure.** Table 3 demonstrates the crucial role of the tree structure in our model’s performance. Removing the CTT component notably decreases performance, especially on multi-class classification datasets like 20News and AMZ, indicating CTT’s importance in capturing document context at the phrase level. Similarly, excluding the DTT component leads to the reduction of the performance, with scores dropping by 4-6% across various datasets, which underscores DTT’s role in capturing inter-word relations while

Models	Hyperpartisan	20News	AMZ	BOOK	ECtHR	Essays
<b>BERT w/ pre-training</b> (Devlin et al., 2019)	91.8	84.7	51.1	58.2	71.7	69.31
<b>ToBERT</b> (Pappagari et al., 2019)	89.5	85.5	54.6	57.3	77.2	72.2
<b>Longformer</b> (Beltagy et al., 2020)	93.7	83.4	56.4	58.5	81.5	74.38
<b>BERT+Random</b> (Park et al., 2022)	89.3	85.0	56.8	59.2	72.8	70.1
<b>BERT+TextRank</b> (Park et al., 2022)	91.2	84.7	56.9	58.9	73.5	70.9
<b>H3-pooler</b> (Lu et al., 2023)	94.2	84.1	57.7	60.5	82.1	-
<b>Ours</b>	<b>95.4</b>	<b>87.0</b>	<b>59.7</b>	<b>62.9</b>	<b>84.9</b>	<b>82.0</b>

Table 2: Comparison to previous methods on the six long document classification datasets with pre-training. Bold: best scores in each column.

Methods		Hyperpartisan	20News	AMZ	BOOK	ECtHR	Essays
<b>Full Model</b>		95.4	87.0	59.7	62.9	84.9	82.0
<b>Tree Structure</b>	Removing CTT	91.3	83.6	54.8	58.8	81.8	78.9
	Removing DTT	90.9	83.1	54.1	58.7	81.7	78.2
	Removing Entire Tree Structure	88.4	78.5	47.3	52.4	77.3	73.7
<b>Removing GAT</b>		88.8	78.9	47.4	50.9	77.8	74.4
<b>Removing the Bidirectional Propagation</b>		87.9	75.9	45.2	48.9	76.1	72.7

Table 3: Ablation experiment results on the six long document classification datasets.

classifying long text. The most significant degradation occurs when the entire tree structure is replaced, especially in the multi-label classification tasks (BOOK, ECtHR, and Essays), indicating that hierarchical document modeling is important for maintaining the structural integrity and contextual coherence of the long documents.

**Effectiveness of the Graph Structure.** As shown in Table 3, removing the GAT module results in performance declines across all datasets, with a more pronounced impact on tasks requiring relational information, such as multi-label classification. This suggests that the GAT module is essential for capturing relationships between different elements within a document, which allows the model to leverage dependencies and interactions that are important for accurate classification. The attention mechanisms provided by the GAT module help focus on important features and connections, which are particularly important for the tasks involving complex relational data.

**Effectiveness of the Bidirectional Information Propagation.** The bidirectional information propagation approach further enhances our model’s effectiveness by fine-tuning feature representations, as reported in Table 3. Removing this module results in significant performance reductions across all datasets, with the most marked impact observed on the BOOK and AMZ datasets. This decline indicates that the bidirectional propagation plays a critical role in polishing the feature representations

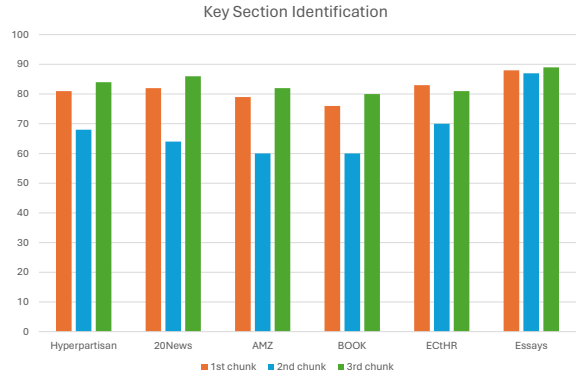


Figure 3: Visualization of the percentage of sentences selected from different chunks across all datasets.

obtained from previous layers, ensuring accurate capture of subtle and important details. By refining these features, the module helps the model better distinguish between different classes, particularly in datasets with high variability and complexity.

### 5.3 Key Section Identification for Long Document

In long document classification, not all sections are equally important. It is crucial to identify which parts of the documents contribute the most significant features for classification tasks. We conduct experiments by splitting each document into three chunks and calculating the number of sentences selected from each chunk. The results, shown in Figure 3, indicate that the first and third chunks consistently contain more important information

530 compared to the middle chunk across all datasets.

531 Specifically, across all datasets except for Essays, 580  
532 approximately 76.6% of sentences are selected 581  
533 from the first chunk, 83% from the third chunk, 582  
534 and only about 62.6% from the middle chunk. In 583  
535 contrast, the Essays dataset, being a questionnaire- 584  
536 style document, distributes important information 585  
537 more evenly, as each question provides different 586  
538 levels of information with equal importance. These 587  
539 observations highlight the importance of effectively 588  
540 capturing content at both the start and end of doc- 589  
541 uments to enhance classification performance in 590  
542 long document tasks. 591

## 543 6 Related Work 592

544 Long documents present unique challenges. As 593  
545 the document length increases, maintaining con- 594  
546 text becomes increasingly difficult, which makes 595  
547 the task substantially more complex compared to 596  
548 short text classification (Liu et al., 2023). Early 597  
549 methods relied on feature extraction techniques, 598  
550 where document length was not a significant is- 599  
551 sue. However, this changed with deep learning 600  
552 approaches using CNNs and RNNs that were im- 601  
553 plemented at different semantic levels, including 602  
554 character, word, and sentence levels (Tang et al., 603  
555 2015; Yang et al., 2016). CNNs often focus on 604  
556 local dependencies (Gao et al., 2018; Liu et al., 605  
557 2018), while RNNs (He et al., 2019; Khandve 606  
558 et al., 2022) are built for long-range dependencies. 607  
559 Transformer-based models, such as BERT (Devlin 608  
560 et al., 2019) and XLNet (Yang et al., 2019), have 609  
561 since come to dominate, but are limited by the 610  
562 number of tokens they can process, which poses 611  
563 a challenge for handling long documents. Meth- 612  
564 ods for modifying transformer architectures to han- 613  
565 dle long documents include recurrent Transform- 614  
566 ers and sparse attention Transformers (Dai et al., 615  
567 2022). Transformer-XL (Dai et al., 2019), a re- 616  
568 current approach, introduced a segment-level re- 617  
569 currence mechanism and a positional encoding 618  
570 scheme, which enabled the model to capture long- 619  
571 term dependencies more effectively. Another ap- 620  
572 proach, ERNIE-Doc (Ding et al., 2021), incorpo- 621  
573 rated a continuous multi-segment attention mecha- 622  
574 nism and entity-aware pre-training to capture com- 623  
575 prehensive contextual information across longer 624  
576 texts. Alternatively, Sparse Transformers (Child 625  
577 et al., 2019) reduced the computational complex- 626  
578 ity of self-attention by selectively focusing on a 627  
579 subset of relevant tokens rather than all tokens in 628

a sequence. Reformer (Kitaev et al., 2020) then 580  
581 improved the efficiency of Transformers by using 582  
583 locality-sensitive hashing for sparse attention and 584  
585 reversible layers to reduce memory usage. Besides 586  
587 the aforementioned approaches, methods such as 588  
589 Longformer (Beltagy et al., 2020) and Big Bird (Za- 590  
591 heer et al., 2020) used a combination of local and 592  
593 global attention mechanisms to reduce computa- 594  
595 tional complexity of standard self-attention. Re- 596  
597 cently, Lu et al. (2023) used SSM to address the 598  
599 computational challenges (quadratic complexity in 600  
601 self-attention) caused by processing long sequences 602  
603 with traditional transformers. They introduced an 604  
605 SSM-pooler model, which incorporates a max pool- 606  
607 ing layer between each SSM block. This design 608  
609 allows the model to automatically extract impor- 610  
611 tant information from nearby inputs at each level 612  
613 and reduce the input length to half of that in the 614  
615 previous layer, which significantly accelerates the 616  
617 speed of both training and inference. They com- 618  
619 pared their model with self-attention-based models 620  
621 and achieved comparable performance while being, 622  
623 on average, 36% more efficient. 624

## 603 7 Conclusion 603

604 In this paper, we address challenges of long docu- 604  
605 ment classification by leveraging a novel graph-tree 605  
606 structure. By integrating syntax trees for sentence 606  
607 encodings and document graphs for comprehensive 607  
608 document encodings, our approach captures both 608  
609 fine-grained syntactic relationships and broader 609  
610 contextual dependencies. Using Tree Transform- 610  
611 ers and GAT ensures accurate modeling of hier- 611  
612 archical relationships within documents. Addi- 612  
613 tionally, our bidirectional information propagation 613  
614 technique enhances the contextual representation, 614  
615 which enables a deeper understanding of content at 615  
616 all hierarchical levels. Notably, our approach not 616  
617 only overcomes the limitations of token constraints 617  
618 but also improves the performance and accuracy 618  
619 of long document classification, making it highly 619  
620 suitable for long document understanding. Poten- 620  
621 tial extensions of this work could involve incorpo- 621  
622 rating external knowledge through the integration 622  
623 of knowledge graphs. By linking document con- 623  
624 tent with relevant external information, the model 624  
625 can further enhance its understanding and context- 625  
626 awareness and it would open new avenues for appli- 626  
627 cations and improve the model’s versatility across 627  
628 various domains. 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

## Limitations

Requiring two forward passes and parsing for the tree-structured transformers increases the time required compared to the other models. This computational overhead should be taken into account when considering the deployment and scalability of the proposed models in practical applications. However, with label-wise attention cutoff values, some sentence and word nodes are pruned, which reduces the computational times significantly and the model takes similar time compared to the BERT-based model for document classification task. Still, with some parallelization in the model implementation, the computational time can be reduced.

## Ethics Statement

We are using the publicly-available datasets, and we do not see any ethics issues in this paper.

## References

Mahtab Ahmed, Muhammad Rifayat Samee, and Robert E. Mercer. 2019a. [You only need attention to traverse trees](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 316–322, Florence, Italy. Association for Computational Linguistics.

Mahtab Ahmed, Muhammad Rifayat Samee, and Robert E Mercer. 2019b. You only need attention to traverse trees. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 316–322.

Muhammad Al-Qurishi. 2022. [Recent advances in long documents classification using deep-learning](#). In *Proceedings of the 5th International Conference on Natural Language and Speech Processing (ICNLSP 2022)*, pages 107–112, Trento, Italy. Association for Computational Linguistics.

Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. Layer normalization. *Advances in NIPS 2016 Deep Learning Symposium*.

David Bamman and Noah A Smith. 2013. New alignment methods for discriminative book summarization. *arXiv preprint arXiv:1305.1319*.

Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. [Longformer: The long-document transformer](#). *Preprint*, arXiv:2004.05150.

Tom Brown, Benjamin Mann, et al. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901.

Ilias Chalkidis, Xiang Dai, Manos Fergadiotis, Prodromos Malakasiotis, and Desmond Elliott. 2022. [An exploration of hierarchical attention transformers for efficient long document classification](#). *arXiv preprint*. 677  
678  
679  
680

Ilias Chalkidis, Manos Fergadiotis, Dimitrios Tsarapat-  
sanis, Nikolaos Aletras, Ion Androutsopoulos, and  
Prodromos Malakasiotis. 2021. [Paragraph-level rationale extraction through regularization: A case study on European court of human rights cases](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 226–241, Online. Association for Computational Linguistics. 681  
682  
683  
684  
685  
686  
687  
688  
689  
690

Rewon Child, Scott Gray, Alec Radford, and Ilya  
Sutskever. 2019. [Generating long sequences with sparse transformers](#). *Preprint*, arXiv:1904.10509. 691  
692  
693

Xiang Dai, Ilias Chalkidis, Sune Darkner, and Desmond  
Elliott. 2022. [Revisiting transformer-based models for long document classification](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 7212–7230, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. 694  
695  
696  
697  
698  
699

Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Car-  
bonell, Quoc Le, and Ruslan Salakhutdinov. 2019. [Transformer-XL: Attentive language models beyond a fixed-length context](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2978–2988, Florence, Italy. Association for Computational Linguistics. 700  
701  
702  
703  
704  
705  
706

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. 707  
708  
709  
710  
711  
712  
713  
714  
715

SiYu Ding, Junyuan Shang, Shuohuan Wang, Yu Sun,  
Hao Tian, Hua Wu, and Haifeng Wang. 2021. [ERNIE-Doc: A retrospective long-document modeling transformer](#). 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 2914–2927, Online. Association for Computational Linguistics. 716  
717  
718  
719  
720  
721  
722  
723  
724

Zican Dong, Tianyi Tang, Lunyi Li, and Wayne Xin  
Zhao. 2023. A survey on long text modeling with  
transformers. *arXiv preprint arXiv:2302.14502*. 725  
726  
727

Daniel Y Fu, Tri Dao, Khaled Kamal Saab, Armin W  
Thomas, Atri Rudra, and Christopher Re. 2023. [Hungry hungry hippos: Towards language modeling with state space models](#). In *The Eleventh International Conference on Learning Representations*. 728  
729  
730  
731  
732

733	Shang Gao, Arvind Ramanathan, and Georgia Tourassi.	Peng Lu, Suyuchen Wang, Mehdi Rezagholizadeh,	787
734	2018. <a href="#">Hierarchical convolutional attention net-</a>	Bang Liu, and Ivan Kobyzev. 2023. <a href="#">Efficient clas-</a>	788
735	<a href="#">works for text classification</a> . In <i>Proceedings of the</i>	<a href="#">sification of long documents via state-space models</a> .	789
736	<i>Third Workshop on Representation Learning for NLP</i> ,	In <i>Proceedings of the 2023 Conference on Empiri-</i>	790
737	pages 11–23, Melbourne, Australia. Association for	<i>cal Methods in Natural Language Processing</i> , pages	791
738	Computational Linguistics.	6559–6565, Singapore. Association for Computa-	792
		tional Linguistics.	793
739	Jun He, Liqun Wang, Liu Liu, Jiao Feng, and Hao Wu.	Agnes Lydia and Sagayaraj Francis. 2019. Ada-	794
740	2019. <a href="#">Long document classification from local word</a>	<a href="#">grad—an optimizer for stochastic gradient descent</a> .	795
741	<a href="#">glimpses via recurrent attention learning</a> . <i>IEEE Ac-</i>	<i>Int. J. Inf. Comput. Sci</i> , 6(5):566–568.	796
742	<i>cess</i> , 7:40707–40718.		
743	Ruining He and Julian McAuley. 2016. <a href="#">Ups and downs:</a>	Christopher Manning, Mihai Surdeanu, John Bauer,	797
744	<a href="#">Modeling the visual evolution of fashion trends with</a>	Jenny Finkel, Steven Bethard, and David McClosky.	798
745	<a href="#">one-class collaborative filtering</a> . In <i>Proceedings of</i>	2014. <a href="#">The Stanford CoreNLP natural language pro-</a>	799
746	<i>the 25th International Conference on World Wide</i>	<a href="#">cessing toolkit</a> . In <i>Proceedings of 52nd Annual Meet-</i>	800
747	<i>Web</i> , WWW ’16, page 507–517, Republic and Can-	<i>ing of the Association for Computational Linguis-</i>	801
748	ton of Geneva, CHE. International World Wide Web	<i>tics: System Demonstrations</i> , pages 55–60, Balti-	802
749	Conferences Steering Committee.	more, Maryland. Association for Computational Lin-	803
		guistics.	804
750	Snehal Ishwar Khandve, Vedangi Kishor Wagh,	Rada Mihalcea and Paul Tarau. 2004. <a href="#">TextRank: Bring-</a>	805
751	Apurva Dinesh Wani, Isha Mandar Joshi, and Ravi-	<a href="#">ing order into text</a> . In <i>Proceedings of the 2004 Con-</i>	806
752	raj Bhuminand Joshi. 2022. <a href="#">Hierarchical neural net-</a>	<i>ference on Empirical Methods in Natural Language</i>	807
753	<a href="#">work approaches for long document classification</a> . In	<i>Processing</i> , pages 404–411, Barcelona, Spain. Asso-	808
754	<i>Proceedings of the 2022 14th International Confer-</i>	ciation for Computational Linguistics.	809
755	<i>ence on Machine Learning and Computing</i> , ICMLC		
756	’22, page 115–119, New York, NY, USA. Association	Raghavendra Pappagari, Piotr Zelasko, Jesús Villalba,	810
757	for Computing Machinery.	Yishay Carmiel, and Najim Dehak. 2019. <a href="#">Hierar-</a>	811
		<a href="#">chical transformers for long document classification</a> .	812
758	Johannes Kiesel, Maria Mestre, Rishabh Shukla, Em-	In <i>2019 IEEE Automatic Speech Recognition and</i>	813
759	manuel Vincent, Payam Adineh, David Corney,	<i>Understanding Workshop (ASRU)</i> , pages 838–844.	814
760	Benno Stein, and Martin Potthast. 2019. <a href="#">SemEval-</a>		
761	<a href="#">2019 task 4: Hyperpartisan news detection</a> . In	Hyunji Park, Yogarshi Vyas, and Kashif Shah. 2022.	815
762	<i>Proceedings of the 13th International Workshop on</i>	<a href="#">Efficient classification of long documents using trans-</a>	816
763	<i>Semantic Evaluation</i> , pages 829–839, Minneapolis,	<a href="#">formers</a> . In <i>Proceedings of the 60th Annual Meet-</i>	817
764	Minnesota, USA. Association for Computational Lin-	<i>ing of the Association for Computational Linguistics</i>	818
765	guistics.	<i>(Volume 2: Short Papers)</i> , pages 702–709, Dublin,	819
		Ireland. Association for Computational Linguistics.	820
766	Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya.	James W Pennebaker and Laura A King. 1999. Lin-	821
767	2020. <a href="#">Reformer: The efficient transformer</a> . In <i>Inter-</i>	<a href="#">guistic styles: language use as an individual differ-</a>	822
768	<i>national Conference on Learning Representations</i> .	<a href="#">ence</a> . <i>Journal of personality and social psychology</i> ,	823
		77(6):1296.	824
769	Ken Lang. 1995. <a href="#">Newsweeder: Learning to filter net-</a>	Duyu Tang, Bing Qin, and Ting Liu. 2015. <a href="#">Document</a>	825
770	<a href="#">news</a> . In Armand Frieditis and Stuart Russell, editors,	<a href="#">modeling with gated recurrent neural network for</a>	826
771	<i>Machine Learning Proceedings 1995</i> , pages 331–339.	<a href="#">sentiment classification</a> . In <i>Proceedings of the 2015</i>	827
772	Morgan Kaufmann, San Francisco (CA).	<i>Conference on Empirical Methods in Natural Lan-</i>	828
		<i>guage Processing</i> , pages 1422–1432, Lisbon, Portu-	829
773	Liu Liu, Kaile Liu, Zhenghai Cong, Jiali Zhao, Yefei	gal. Association for Computational Linguistics.	830
774	Ji, and Jun He. 2018. <a href="#">Long length document classi-</a>	Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen,	831
775	<a href="#">fication by local convolutional feature aggregation</a> .	Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang,	832
776	<i>Algorithms</i> , 11(8).	Sebastian Ruder, and Donald Metzler. 2020. Long	833
		range arena: A benchmark for efficient transformers.	834
777	Tengfei Liu, Yongli Hu, Boyue Wang, Yanfeng Sun, Jun-	In <i>International Conference on Learning Representa-</i>	835
778	bin Gao, and Baocai Yin. 2023. <a href="#">Hierarchical graph</a>	<i>tions</i> .	836
779	<a href="#">convolutional networks for structured long document</a>		
780	<a href="#">classification</a> . <i>IEEE Transactions on Neural Net-</i>	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	837
781	<i>works and Learning Systems</i> , 34(10):8071–8085.	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	838
		Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	839
782	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	Bhosale, et al. 2023. LLaMa 2: Open founda-	840
783	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	tion and fine-tuned chat models. <i>arXiv preprint</i>	841
784	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	<i>arXiv:2307.09288</i> .	842
785	Roberta: A robustly optimized BERT pretraining		
786	approach. <i>arXiv preprint arXiv:1907.11692</i> .		

843 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob  
844 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
845 Kaiser, and Illia Polosukhin. 2017. Attention is all  
846 you need. *Advances in Neural Information Process-*  
847 *ing Systems*, 30.

848 Petar Veličković, Guillem Cucurull, Arantxa Casanova,  
849 Adriana Romero, Pietro Liò, and Yoshua Bengio.  
850 2018. [Graph attention networks](#). In *International*  
851 *Conference on Learning Representations*.

852 Danqing Wang, Pengfei Liu, Yining Zheng, Xipeng Qiu,  
853 and Xuanjing Huang. 2020. [Heterogeneous graph](#)  
854 [neural networks for extractive document summariza-](#)  
855 [tion](#). In *Proceedings of the 58th Annual Meeting of*  
856 *the Association for Computational Linguistics*, pages  
857 6209–6219, Online. Association for Computational  
858 Linguistics.

859 Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Car-  
860 bonell, Ruslan Salakhutdinov, and Quoc V. Le. 2019.  
861 XLNet: generalized autoregressive pretraining for  
862 language understanding. In *Proceedings of the 33rd*  
863 *International Conference on Neural Information Pro-*  
864 *cessing Systems*, Red Hook, NY, USA. Curran Asso-  
865 ciates Inc.

866 Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He,  
867 Alex Smola, and Eduard Hovy. 2016. [Hierarchical](#)  
868 [attention networks for document classification](#). In  
869 *Proceedings of the 2016 Conference of the North*  
870 *American Chapter of the Association for Computa-*  
871 *tional Linguistics: Human Language Technologies*,  
872 pages 1480–1489, San Diego, California. Associa-  
873 tion for Computational Linguistics.

874 Manzil Zaheer, Guru Guruganesh, Avinava Dubey,  
875 Joshua Ainslie, Chris Alberti, Santiago Ontanon,  
876 Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang,  
877 and Amr Ahmed. 2020. Big Bird: transformers for  
878 longer sequences. In *Proceedings of the 34th Interna-*  
879 *tional Conference on Neural Information Processing*  
880 *Systems, NIPS '20*, Red Hook, NY, USA. Curran  
881 Associates Inc.