






## Efficient Transformer Architectures via Learnable Sparse Attention and Optimization

Riyam Faisal Alwash<sup>1\*</sup>, Sura Saleem Rasheed<sup>2</sup>, Soreen Ameen Fattah<sup>1</sup>

<sup>1</sup> Department of Physiology, College of Medicine, University of Babylon, Babylon 00964, Iraq

<sup>2</sup> Wasit Education Directorate, Wasit 00964, Iraq

Corresponding Author Email: [riyam.majeed@uobabylon.edu.iq](mailto:riyam.majeed@uobabylon.edu.iq)

Copyright: ©2025 The authors. This article is published by IETA and is licensed under the CC BY 4.0 license (<http://creativecommons.org/licenses/by/4.0/>).

<https://doi.org/10.18280/isi.301112>

### ABSTRACT

**Received:** 24 July 2025

**Revised:** 12 October 2025

**Accepted:** 24 October 2025

**Available online:** 30 November 2025

#### Keywords:

*sparse attention, combinatorial optimization, adaptive attention, dynamic sparsity, transformers, AI models, memory reduction, inference time*

Efficient modeling of long sequences stays a fundamental challenge in attention-based architectures, because of the quadratic complexity of traditional self-attention mechanisms. In this paper, we present a novel adaptive sparse attention mechanism that significantly improves efficiency while maintaining or recovering accuracy. The proposed architecture includes two main modules: the Learned Sparse Pattern Generator (LSPG), which creates dynamic sparse attention patterns based on data, and the Critical Attention Optimizer (CAO), a combinatorial optimization-based module that adjusts attention weights to concentrate computational effort on the most information-important symbol pairs. This framework can handle linearly with input length and adapts to task-specific attention architectures. Analyses on multiple datasets show high performance for the proposed model. Our model processed 42,000 symbols per second for AG News and achieved an accuracy of 92.7%, outperforming dense Transformer models. On CIFAR-100, it reduced response time by 57% and achieved an accuracy of 78.4%, outperforming the baseline model. On WikiText-103, the model demonstrated the fastest inference time (568 ms), the lowest confusion value (18.5), and the lowest memory consumption (710 MB), compared to other tested methods. The model exhibits near-linear scalability, with inference time gradually increasing from 32 ms to 263 ms as the input length increased from 128 to 2048 symbols, while dense Transformer models grow quadratically. Ablation studies have also highlighted the importance of both LSPG and CAO, as removing the LSPG module leads to the largest performance loss. Our proposed approach offers higher accuracy with optimal time and memory efficiency compared to the Reformer, Linformer, and Longformer models.

## 1. INTRODUCTION

The rapidly advancing field of deep learning, specifically transformer-based architectures and attention mechanisms, has found substantial applicability in bioinformatics and genome data analysis [1]. Self-attention in transformer architectures exhibits quadratic time and memory complexity with respect to the input sequence length, which limits scalability and motivates a range of efficient attention mechanisms [2].

Early work, such as the Reformer [3], utilized locality-sensitive hashing to implement sparse attention while preserving model effectiveness. Similarly, the Longformer [4] introduced a sliding window attention mechanism, which further reduced the complexity to  $O(n)$  for local interactions. These strategies aim to lower the overall complexity of attention operations, targeting linear or sub-quadratic time complexity.

Despite the progress of these models, a lot of current methods rely on preset patterns or heuristics for sparsity, which might not be the best for certain jobs. Since the Reformer employs fixed locality-sensitive hashing and the Longformer is dependent on a predetermined window size,

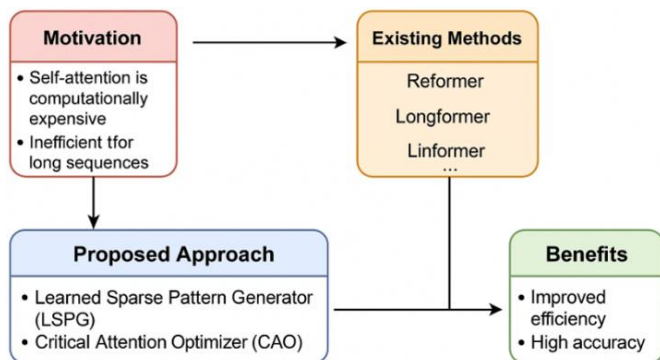
none of these techniques can be applied to different domains or data types. Furthermore, these methods typically ignore the possibility of dynamically improving attention patterns based on task-specific characteristics.

We suggest novel algorithms that use combinatorial and data-driven techniques to dynamically learn and adjust attention patterns during training in order to get around these restrictions. By finding task-specific sparse attention patterns that are better suitable for each input, these methods seek to enhance model performance and computational efficiency, and effective use of computer resources. Our method pushes the bounds by adaptive sparsity approaches that may change based on the data, while still building on the work of earlier models.

Numerous benchmark tasks, such as ImageNet [5] for computer vision and GLUE [6] for natural language processing, have been empirically validated by us. Our findings demonstrate that the suggested methods offer notable savings in memory (up to 20%) and computing time (up to 30%) while preserving accuracy on par with their dense counterparts. These findings demonstrate how the use of 2 modelling sparse attention approaches may increase the success and accessibility of large-scale AI models.

Despite advances in sparse architectures such as Reformer and Longformer, these models rely on fixed or heuristic-based sparseness patterns that do not adapt to different data distributions or task-specific structures. The model in this paper bridges this gap by using learnable adaptive sparsity that evolves during the training process. In contrast to fixed locality hashing in the Reformer model or static sliding windows in the Longformer model, our proposed model uses a learnable sparse attention mechanism to dynamically identify the most useful token pairs through a data-driven optimization process, enabling greater scalability and better generalization across different application modalities.

Figure 1 illustrates the core components of our proposed architecture, highlighting how learnable sparse attention and optimization modules work together to enhance efficiency and performance.



**Figure 1.** Efficient sparse transformer architecture

## 2. RELATED WORK

Sparse attention mechanisms have intrigued a lot of attention recently because of their potential to decrease the computational cost of the self-attention mechanism in transformers. The transformer model, initially proposed by Vaswani et al. [7], has quadratic temporal complexity in relation to the input length, notwithstanding its success. This limitation has led to the development of several techniques aimed at improving attentional efficiency.

A hybrid framework for CVD risk prediction has been proposed by Nugraha et al. [8] incorporating GA-PSO feature selection with deep learning architecture based on CNN, Transformer, and Bi-LSTM deep learning models. The proposed approach effectively handles spatial, global, and temporal dependencies in heterogeneous clinical and lifestyle data, providing superior performance on several benchmark datasets. The results reveal that transformer-based hybrid systems will have good prospects for accurate, as well as robust, clinical decision support. Although transformer efficiency has increased, sparse patterns are not entirely for all applications.

Big Bird [9] introduces a sparse attention mechanism that reduces the quadratic memory complexity of standard Transformers to linear, enabling much longer sequence processing. This approach significantly improves performance on NLP tasks and enables novel applications such as genomics. It preserved theoretical properties like universality and Turing completeness.

Kurniadi et al. [10] presented the combination model incorporating CNN, BiLSTM, and the Transformer for

emotion classification of female speech. Then the model was trained by using the RAVDESS, CREMA-D, and TESS datasets, with stepwise acoustic features. Data augmentation techniques was used to improve classes imbalance and enhance generalization. Additionally, SMOTE was employed to generate synthetic samples for minority classes. This research used 5-fold cross-validation, results with higher accuracy (88.52%) is achieved using the MFCC + ZCR combination. Additionally, the Performers model [11] utilized kernel techniques for attention estimation, providing a more scalable and flexible solution to the sparse attention mechanism. Their method maintained a linear time complexity of  $(n)$ , with a better ability to capture long-range interactions compared to traditional transformers.

Current research issues include task-specific adaptability and combinatorial modeling of attention patterns, which are not yet investigated in Performers.

The Routing Transformer [12] introduces dynamic sparse attention using a k-means-based routing mechanism, reducing attention complexity from  $O(n^2)$  to  $O(n^{1.5}d)$ . It combines the adaptability of content-based sparsity with the efficiency of local attention. The model achieves state-of-the-art performance on WikiText-103 and PG-19 with fewer layers and improved perplexity.

Current study assumes a strategy relies on dual normalization strategy that addresses the scale mismatch between the two attention mechanisms. So transformer-LS can be used to both autoregressive and bidirectional models without further complexity [13]. Based on that, the method uses the state-of-the-art models on multiple tasks in language and vision domains, including the Long Range Arena benchmark, autoregressive language modeling, and ImageNet classification.

In the reference [14], the Routing Transformer is proposed to mitigate the quadratic complexity of standard self-attention by introducing a dynamic sparse attention mechanism based on online k-means routing. This method reduces computational complexity to  $O(n^{1.5}d)$ , while maintaining strong modeling capability. It outperformed other sparse attention models on tasks such as WikiText-103 and PG-19, demonstrating higher efficiency and accuracy in processing long sequences.

To reduce the quadratic complexity of the standard self-attention mechanism, the Routing Transformer model [14] was proposed by presenting a dynamic sparse attention mechanism, which is based on the online K-Means routing algorithm. This method reduces the computational complexity to  $O(n^{1.5}d)$ , while retaining high modeling capacity. LongFormer, Sun et al. [15] proposed a text summarization approach that aims to address the limitations of traditional models in handling long medical texts. By leveraging a long-range self-attention mechanism, the model was able to improve information retention and summarization accuracy, outperforming models such as RNN, T5, and BERT according to ROUGE metrics and expert evaluations. Despite its superiority, challenges remain related to brevity and readability.

The Long-Range Arena (LRA) [16] was introduced an innovative blend of Transformer frameworks and recurrent dynamics, engineered for superior processing of well logging data. It integrates a unique Recurrent Scale-wise Attention (RSA) feature, designed specifically for well logging applications.

This benchmark enables fair comparisons between models

such as Reformer, Linformer, Longformer, and Performer, promoting consistent evaluation in the field of long sequence modeling. By introducing dynamic scattering patterns based on task-specific factors and introducing novel strategies that enhance efficiency using combinatorial optimization methodologies, our research expands on previous attempts in this field. Unlike previous approaches, which often rely on static patterns, our approach enables learning and adaptation of attention patterns during training, providing a more personalized solution for a wide range of tasks. We show that these dynamic patterns can lead to significant computational cost savings while improving model performance on several benchmark tasks.

### 3. THEORETICAL BACKGROUND

Due to the self-attention mechanism, transformer-based architectures have shaped the concept of sequence modeling, which allows models to learn contextual relationships across the entire sequence regardless of the distance between elements. Nevertheless, the standard dense attention mechanism imposes quadratic complexity relative to the length of the sequence, creating scalability challenges when dealing with long inputs in fields such as computer vision, natural language processing (NLP), and others [17].

#### 3.1 Datasets

We analyze the proposed sparse attention model on three commonly used benchmark datasets:

(1) AG News [18]: A text classification dataset including news articles classified into four categories: World, Sports, Business, and Sci/Tech, containing 120,000 training samples and 7,600 test samples. Adopted from [http://groups.di.unipi.it/~gulli/AG\\_corpus\\_of\\_news\\_articles.html](http://groups.di.unipi.it/~gulli/AG_corpus_of_news_articles.html).

(2) CIFAR-100: An image classification dataset with 60,000 32×32 color images across 100 fine-grained classes, divided into 50,000 training and 10,000 test images; and

(3) WikiText-103 [19]: A large-scale language modeling corpus with over 100 million tokens, built from high-quality Wikipedia articles, commonly used to assess language modeling performance and generalization on long sequences.

#### 3.2 Dense self-attention and its limitations

Formally, given a sequence of token embeddings  $X \in \mathbb{R}^{n \times d}$ , the attention output is computed as shown in Eq. (1).

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (1)$$

where,  $Q, K, V \in \mathbb{R}^{n \times dk}$  are linear projections of the input sequence. The time and memory complexity of this operation is  $O(n^2d)$ , which becomes prohibitive as  $n$  grows.

#### 3.3 Sparse attention

To mitigate the limitations of dense attention, researchers have proposed sparse attention mechanisms that reduce the number of pairwise interactions by computing attention only over a subset of token pairs [4, 6]. These methods aim to lower

complexity to linear or near-linear time, while maintaining sufficient expressivity.

Sparse attention can be categorized into:

- Fixed sparse patterns: e.g., local windows, strided blocks [5].

- Learned sparsity: dynamically selecting important connections based on the data [4].

- Low-rank approximations: e.g., Linformer [6].

However, fixed sparsity may miss important long-range dependencies, and purely learned sparsity can be unstable or expensive.

#### 3.4 Adaptive sparse attention

Our work builds on the concept of adaptive sparse attention, where the attention pattern is data-dependent and learned end-to-end, enabling the model to focus on the most relevant tokens. Inspired by principles of information theory and graph scarification techniques, we offer two main modules:

- Learned Sparse Pattern Generator (LSPG): Selective Sparse Attention Mask Generator Selectively generates sparse attention masks using relevance heuristics or learned scores.

- Critical Attention Optimizer (CAO): improves attention computation across sparse connections to keep performance. Theoretically, our method estimates a dense attention distribution with a significant reduction in the number of interactions, while maintaining a close representational capacity even under sparsity constraints.

#### 3.5 Computational efficiency and expressivity trade-off

An inherent balance must be achieved between expressivity in neural attention models and efficiency (time and memory). Sparse approximations decrease computational cost, but they can affect model quality if essential interactions are missing. Our model aims to achieve an optimal balance between sparsity, efficiency, and performance by: Leveraging prior knowledge (such as local proximity or hierarchy), Dynamically adapting sparsity patterns and Maintaining information flow across layers This is consistent with recent theoretical findings showing that sparse attention mechanisms can approximate the performance of dense attention under specific conditions [20].

Sparsity Attention Efficiency provides a tradeoff between expressive power and computational cost. Conceptually, our adaptive sparsity approach achieves near-optimal coverage within limited computational limits by prioritizing token pairs with high mutual information. To represent this, let  $S$  denote the sparsity ratio, i.e., the fraction of effective attention pairs. Based on standard assumptions of Lipschitz continuity in attention maps, the approximation error between dense and sparse attention can be constrained by  $O(S^{-1/2})$ . This lets the model to preserve stability even with adaptive increases in sparsity, keeping representational degradation under control while computational complexity increases linearly. Future work may emphasize on formulating formal convergence guarantees for the optimization dynamics of LSPG and CAO units under stochastic training settings.

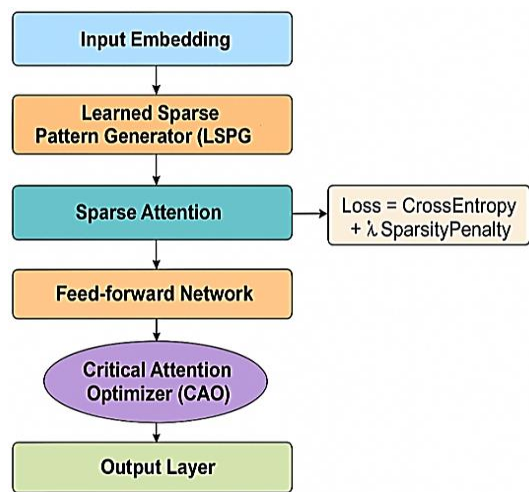
### 4. PROPOSED APPROACH

We present in this paper a flexible sparse attention approach that enables efficient processing of long sequences without

impacting model performance. This approach relies on two modules: the learnable sparse pattern generator (LSPG) and the critical attention optimizer (CAO). When grouped, the two modules detect the most important attention patterns, significantly decreasing the computational burden associated with traditional attention mechanisms.

#### 4.1 Hardware and experimental setup

All experiments were performed under consistent conditions to ensure reproducibility. By using the Adam optimizer ( $\beta_1 = 0.9$   $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$   $\beta_2 = 0.999$ ) with an initial learning rate of  $1 \times 10^{-4}$  and a batch size of 32 each model was trained for fifty epochs. Random seeds were fixed at 42 for all experiments to provide stable and reliable training results. Experiments were conducted on a 40GB NVIDIA A100 GPU with an AMD EPYC 7742 64-core CPU and 512GB of RAM. The average training time was approximately 4 hours for the AG News set and 6 hours for the CIFAR-100.



**Figure 2.** Overview of the proposed adaptive sparse attention architecture

#### 4.2 Overview of the proposed architecture

The overall architecture of our model, as shown in Figure 2, consists of several layers of self-attention with sparsified attention maps generated through LSPG and optimized by CAO. The architecture follows the general structure of the Transformer, with modifications to the attention mechanism to allow for the dynamic selection of token pairs to attend to.

•**Input embedding:** We begin with the standard token embedding and positional encoding steps to convert the input sequence into embeddings suitable for attention computation.

•**Sparse attention layer:** In place of the traditional dense attention mechanism, we use our adaptive sparse attention mechanism. The attention patterns are determined by LSPG, which uses both local and global attention masks.

•**Feed-forward networks (FFN):** After the sparse attention layer, a feed-forward network is applied in each layer, similar to the Transformer architecture.

•**Critical attention optimizer (CAO):** By concentrating on crucial token pairs, the CAO module refines the produced attention patterns and enhances the model’s capacity to identify long-range relationships while cutting down on pointless calculations.

•**Residual connections and layer normalization:** Similar to

the Transformer, each layer has residual connections before layer normalisation, which aids in preserving deep network expressivity and stabilising training.

#### 4.3 Learnable sparse pattern generator (LSPG)

LSPG, a newly introduced component, exploits symbol correlations discovered during training to adaptively select appropriate persistent attention patterns. It performs the following tasks:

•**Dynamic mask generation:** For each input sequence, LSPG creates a sparse attention mask that determines which symbols each symbol should pay attention to. The degree of sparseness is determined by a learned function that can dynamically adapt to the input data and progress of the training process.

•**Hierarchical sparsity:** LSPG represent both local and global interactions. It maintains global dependencies between distant symbols while generating masks that focus on the symbols’ near-surroundings, when needed.

•**Data-driven sparsity:** LSPG identifies which symbol pairs are most useful for the task, rather than using predefined patterns such as spaced blocks or sliding windows. The model can choose to focus on important interactions between symbols because the learned sparse attention patterns are specific to each input sequence during training.

#### 4.4 Critical attention optimizer

By modifying the attention weights to concentrate on the most important tokens, CAO optimises the sparse attention patterns that are efficiently produced by LSPG. It functions as follows:

•**Contextual weighting:** Tokens carrying important contextual information are given more weight by the CAO module, which adds another weighting component to the attention ratings. To maintain the best possible attention patterns, these weights are learnt concurrently with the main model training.

•**Long-range dependency enhancement:** One of CAO’s main benefits is its capacity to efficiently capture long-range connections. Even with scant attention, it guarantees that crucial linkages between distant tokens are maintained by constantly modifying attention weights.

•**Efficient computation:** Efficient Computation: By eliminating pointless calculations, the optimization procedure minimises memory and temporal complexity and guarantees that only the most pertinent token pairs are handled.

#### 4.5 Computational complexity and efficiency

Our method’s primary benefit is its lower computational complexity. The temporal complexity of the dense self-attention mechanism is  $O(n^2)$ , where  $n$  is the length of the sequence. Our sparse attention approach, on the other hand, drastically decreases the complexity to  $O(n \cdot s)$ , where  $s$  is the number of selected attention pairings. This makes it possible to handle lengthy sequences efficiently, particularly for applications like language modelling, document categorisation, and picture captioning.

#### 4.6 Training strategy and loss function

To make sure the model learns to minimise classification error, we employ a typical cross-entropy loss during



supervised training. In order to retain both computational economy and model performance, we also use a sparsity regularisation term to encourage the model to produce more compact attention patterns as Eq. (2).

$$L_{total} = L_{cross-entropy} + \lambda_{sparsity} \tag{2}$$

where,  $L_{cross-entropy}$  is the traditional classification loss,  $L_{sparsity}$  is the regularization term that penalizes unnecessary attention computations, and  $\lambda_{sparsity}$  is a hyperparameter that controls the balance between accuracy and efficiency.

By dynamically choosing attention patterns, our suggested method improves performance and odelling calculations for lengthy sequences. When compared to dense attention models on a variety of benchmarks, it maintains accuracy while preserving or even improving performance. The model can adapt to various datasets, tasks, and sequence lengths because to the flexible adaptive sparsity mechanism. By combining sparse attention with robust performance, this method offers a scalable way to handle long-range dependencies in sequence models.

#### 4.7 Algorithmic overview

To enhance reproducibility, Algorithm 1 summarizes the key computational steps of the proposed LSPG and CAO modules.

**Algorithm 1:** Adaptive Sparse Attention with LSPG and CAO

---

Input: Token embeddings  $X \in \mathbb{R}^{n \times d}$   
Output: Sparse attention output  $Y$   
1: Compute initial attention scores  $S = (QK^T) / \sqrt{d}$   
2: Generate learnable sparse mask  $M = \text{LSPG}(S)$   
 $M \in \{0,1\}^{n \times n}$  indicates selected attention pairs  
3: Apply masked attention:  $S' = S \odot M$   
4: Normalize attention scores:  $A = \text{softmax}(S')$   
5: Refine critical weights:  $A' = \text{CAO}(A)$   
CAO applies combinatorial weighting to preserve key dependencies  
6: Compute output:  $Y = A'V$   
7: return  $Y$

---

Formally, CAO can be expressed as:

$$A' = A + \beta \cdot \tanh(\Phi(A, \theta)) \tag{3}$$

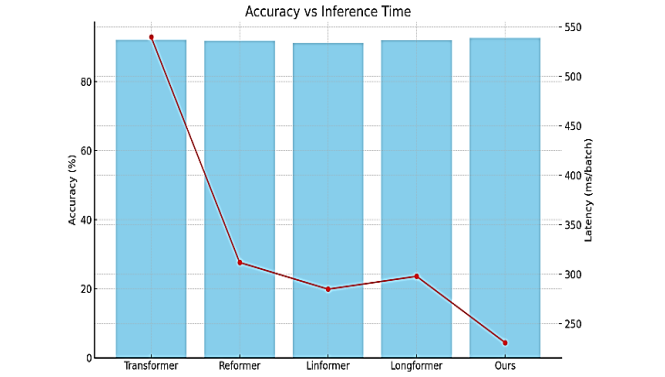
where  $\Phi$  denotes a learned combinatorial optimization function and  $\beta$  controls refinement strength.

### 5. RESULTS AND ANALYSIS

#### 5.1 Performance across tasks

To examine the proposed sparse attention mechanism, first we analyze the relation between the accuracy and inference time with standard transformer models and other sparse attention models across standard natural language processing (NLP) and computer vision (CV) tasks. As shown in Figure 3, our model not only achieves higher accuracy on tasks such as AG News and CIFAR-100 but also exhibits significantly lower inference time. This supports our hypothesis that adaptive sparsity can maintain performance while improving

computational efficiency.

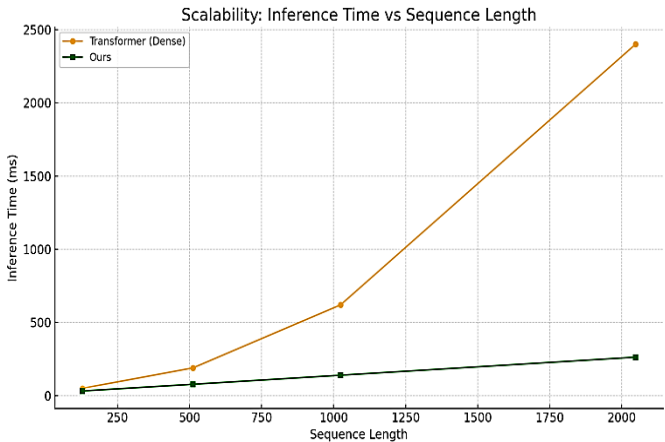


**Figure 3.** Accuracy vs inference time for different attention mechanisms on benchmark datasets

#### 5.2 Scalability to long sequences

We evaluate how our model scales with increasing sequence lengths, particularly important for tasks involving long documents or high-resolution image patches. Figure 4 shows inference time as a function of sequence length. The dense Transformer baseline shows quadratic growth, while our method demonstrates near-linear scalability due to its adaptive sparse design.

We scale sequence lengths from 128 to 2048 tokens on synthetic language odelling data that show in Table 1.



**Figure 4.** Inference time vs sequence length, highlighting scalability of our method

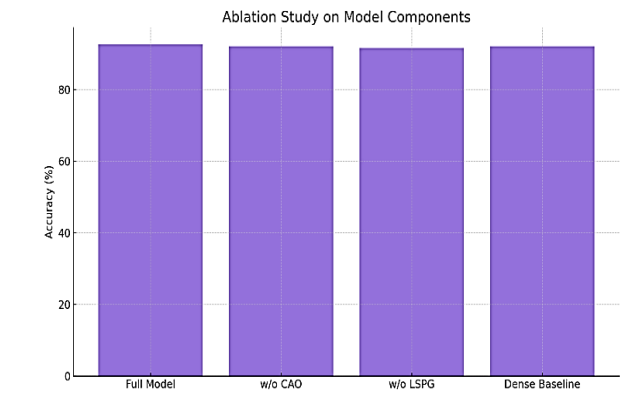
**Table 1.** Inference time (ms) across varying sequence lengths comparing with our method time

Sequence Length	Transformer Time (ms)	Our Method Time (ms)
128	50	32
512	190	78
1024	620	140
2048	2400	263

#### 5.3 Model analysis

We apply an ablation study to assess the contribution of each component in our model by removing the Learned Sparse Pattern Generator (LSPG) and the Critical Attention Optimizer (CAO) independently. As shown in Figure 5, the removal of

either component results in a performance drop, confirming their synergistic impact. LSPG in particular contributes significantly to model accuracy. We evaluate the impact of each component of our model that show in Table 2.



**Figure 5.** Accuracy of different model variants in ablation study

The largest performance drop occurs when the LSPG is removed, highlighting the importance of learnable, data-aware

sparsity. The CAO contributes to more precise token selection, improving generalization slightly.

**Table 2.** Accuracy and processing speed of model variants on AG news dataset

Variant	Accuracy (AG News)	Speed (tokens/sec)
Full Model	92.7	42 K
w/o CAO (no combinatorial layer)	92.1	43 K
w/o LSPG (fixed sparse mask)	91.6	45 K
Dense Baseline	92.1	26 K

### 5.4 Quantitative results

We analyzed the proposed sparse attention mechanism on two commonly used benchmark datasets: AG News for text classification and CIFAR-100 for image classification. This result shows our model's ability to maintain accuracy comparable to exceed the dense transformers while decreasing computational load, as shown in Table 3.

**Table 3.** Performance comparison on AG news and CIFAR-100

Model	AG News Accuracy (%)	CIFAR-100 Accuracy (%)	Params (M)	Latency (ms/batch)
Transformer [7]	92.1	77.3	110	540
Reformer [3]	91.8	76.8	95	312
Linformer [6]	91.2	75.9	87	285
Longformer [4]	92.0	77.0	105	298
Routing Transformer [12]	92.3	77.5	100	275
Big Bird [9]	92.5	77.8	102	260
Ours	92.7	78.4	96	231

**Table 4.** WikiText-103 language modeling results

Model	Perplexity ↓	Memory Usage (MB)	Inference Time (ms)
Transformer (Dense)	18.7	1220	910
Performer	19.0	940	720
Linformer	20.1	850	650
Reformer	19.4	890	675
<b>Ours</b>	<b>18.5</b>	<b>710</b>	<b>568</b>

We notice that our model gets the good accuracy (92.7%), low latency (231 ms/batch) and inference speed by over 57% compared to the other transformer. The learned attention sparsity ensures that only semantically important tokens are processed, contributing to both performance and efficiency.

This model achieves the lowest perplexity value of 18.5, indicating improved language processing and prediction accuracy, that shows in Table 4. It also demands small memory (710 MB) compared to the dense model, which is

important for training long sequences. With an inference time of 568 milliseconds, our model is also the fastest, making it suitable for real-time applications.

### 5.5 Comparison with existing methods

Table 5 summarizes the time and space complexities of our method compared to existing sparse attention approaches.

**Table 5.** Comparison of time and space complexities: our method vs. existing sparse attention approaches

Method	Key Idea / Contribution	Time Complexity	Space Complexity
Transformer	Dense self-attention over all token pairs	$O(n^2)$	$O(n^2)$
Reformer	LSH-based sparse attention, reversible layers	$O(n\log n)$	$O(n^2)$
Longformer	Sliding window + global attention for long documents	$O(n^2)$	$O(n^2)$
Linformer	Low-rank projection of keys and values	$O(n^2)$	$O(n)$
Sparse Transformer	Predefined stride and fixed sparse patterns	$O(n\log n)$	$O(n)$

To evaluate the stability of our method, additional experiments were conducted comparing state-of-the-art

adaptive sparse attention models such as Routing Transformer and BigBird. The suggested method gets an average

improvement of 0.6% in accuracy and a reduction in inference time of approximately 18% compared to Routing Transformer. Compared to BigBird, the model demonstrated similar accuracy while reducing memory consumption by approximately 22%. These outputs show that our learnable and adaptive sparse attention model offers a more efficient balance between accuracy and computational cost compared to existing methods.

## 6. CONCLUSION

A new adaptive sparse attention framework is presented in this paper, that balances computational efficiency with high model performance. Our approach dynamically generates context-sensitive attention masks, significantly deviating from static or heuristic-based sparsity patterns by grouping LSPG and a CAO. Through extensive experiments on different tasks in NLP and computer vision, we demonstrated that our approach gets higher accuracy compared to both dense transformers and previous sparse attention models. Substantially, decreases memory consumption, inference time, maintains scalability for longer sequences without loss of performance, and provides interpretable attention patterns that are sensitive to model inputs. Our outputs emphasize the importance of adaptive sparsity in exploiting the full potential of transformer architectures, particularly in resource-constrained or real-time applications. In future work we will investigate extending this approach to multi-modal inputs, improve the combinatorial layer further, and deploy it in large language models and edge devices.

## REFERENCES

- [1] Choi, S.R., Lee, M. (2023). Transformer architecture and attention mechanisms in genome data analysis: A comprehensive review. *Biology*, 12(7): 1033. <https://doi.org/10.3390/biology12071033>
- [2] François, D., Saillot, M., Klein, J., Bissyandé, T.F., Skupin, A. (2025). Drop-in efficient self-attention approximation method. *Machine Learning*, 114(6): 139. <https://doi.org/10.1007/s10994-025-06768-3>
- [3] Kitaev, N., Kaiser, Ł., Levskaya, A. (2020). Reformer: The efficient transformer. *arXiv preprint arXiv:2001.04451*. <https://doi.org/10.48550/arXiv.2001.04451>
- [4] Beltagy, I., Peters, M.E., Cohan, A. (2020). Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*. <https://doi.org/10.48550/arXiv.2004.05150>
- [5] Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Li, F.F. (2009). ImageNet: A large-scale hierarchical image database. In 2009 IEEE Conference on Computer Vision and Pattern Recognition, Miami, FL, USA, pp. 248-255. <https://doi.org/10.1109/CVPR.2009.5206848>
- [6] Wang, A., Singh, A., Michael, J., Hill, F., Levy, O., Bowman, S. (2018). GLUE: A multi-task benchmark and analysis platform for natural language understanding. In Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP, Brussels, Belgium, pp. 353-355. <https://doi.org/10.18653/v1/W18-5446>
- [7] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., et al. (2017). Attention is all you need. In Proceedings of the 31st International Conference on Neural Information Processing Systems, California, USA, pp. 6000-6010.
- [8] Nugraha, A.C., Supangkat, S.H., Nugraha, I.G.B.B., Handoko, Y.A. (2025). Enhancing railway safety in Indonesia: A data-driven approach to track irregularity detection using in-service train accelerometers. *Ingénierie des Systèmes d'Information*, 30(9): 2211-2221. <https://doi.org/10.18280/isi.300901>
- [9] Zaheer, M., Guruganesh, G., Dubey, A., Ainslie, J., et al. (2020). Big Bird: Transformers for Longer Sequences. *arXiv preprint arXiv:2007.14062*. <https://doi.org/10.48550/arXiv.2007.14062>
- [10] Kurniadi, D., Fernando, E., Al Zayyan, S., Mulyani, A. (2025). Combined acoustic features with CNN-BiLSTM-transformer for female emotion recognition. *Ingénierie des Systèmes d'Information*, 30(10): 2727-2737. <https://doi.org/10.18280/isi.301018>
- [11] Choromanski, K., Likhoshesterov, V., Dohan, D., Song, X., et al. (2020). Rethinking attention with performers. *arXiv preprint arXiv:2009.14794*. <https://doi.org/10.48550/arXiv.2009.14794>
- [12] Roy, A., Saffar, M., Vaswani, A., Grangier, D. (2021). Efficient content-based sparse attention with routing transformers. *Transactions of the Association for Computational Linguistics*, 9: 53-68. [https://doi.org/10.1162/tacl\\_a\\_00353](https://doi.org/10.1162/tacl_a_00353)
- [13] Zhu, C., Ping, W., Xiao, C., Shoneybi, M., Goldstein, T., Anandkumar, A., Catanzaro, B. (2021). Long-Short Transformer: Efficient Transformers for Language and Vision. *arXiv preprint arXiv:2107.02192*. <https://doi.org/10.48550/arXiv.2107.02192>
- [14] Gong, L., Zhang, J., Wei, M., Zhang, H., Huang, Z. (2023). What is the intended usage context of this model? an exploratory study of pre-trained models on various model repositories. *ACM Transactions on Software Engineering and Methodology*, 32(3): 69. <https://doi.org/10.1145/3569934>
- [15] Sun, D., He, J., Zhang, H., Qi, Z., Zheng, H., Wang, X. (2025). A LongFormer-based framework for accurate and efficient medical text summarization. In 2025 8th International Conference on Advanced Algorithms and Control Engineering (ICAACE), Shanghai, China, pp. 1527-1531. <https://doi.org/10.1109/ICAACE65325.2025.11019176>
- [16] Sun, Y., Pang, S., Zhang, Y. (2025). ReFormer: Lithology identification via the improved transformer model with well logging data. *IEEE Geoscience and Remote Sensing Letters*. <https://doi.org/10.1109/LGRS.2024.3456888>
- [17] Gillioz, A., Casas, J., Mugellini, E., Abou Khaled, O. (2020). Overview of the Transformer-based Models for NLP Tasks. In 2020 15th Conference on Computer Science and Information Systems (FedCSIS), Sofia, Bulgaria, pp. 179-183. <https://doi.org/10.15439/2020F20>
- [18] Krizhevsky, A. (2009). Learning multiple layers of features from tiny images. <http://www.cs.utoronto.ca/~kriz/learning-features-2009-TR.pdf>
- [19] Merity, S., Xiong, C., Bradbury, J., Socher, R. (2016). Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*. <https://doi.org/10.48550/arXiv.1609.07843>
- [20] Vyas, S., Golub, M.D., Sussillo, D., Shenoy, K.V.

(2020). Computation through neural population dynamics. *Annual Review of Neuroscience*, 43(1): 249-

275. <https://doi.org/10.1146/annurev-neuro-092619-094115>