# Saliency-driven Dynamic Token Pruning for Large Language Models

**Yao Tao, Yehui Tang, Yun Wang, Mingjian Zhu,**
**Hailin Hu**\*, **Yunhe Wang**\*
Huawei Noah's Ark Lab.
{hailin.hu,yunhe.wang}@huawei.com

## Abstract

Despite the recent success of large language models (LLMs), LLMs are particularly challenging in long-sequence inference scenarios due to the quadratic computational complexity of the attention mechanism. Inspired by the interpretability theory of feature attribution in neural network models, we observe that not all tokens have the same contribution. Based on this observation, we propose a novel token pruning framework, namely Saliency-driven Dynamic Token Pruning (SDTP), to gradually and dynamically prune redundant tokens based on the input context. Specifically, a lightweight saliency-driven prediction module is designed to estimate the importance score of each token with its hidden state, which is added to different layers of the LLM to hierarchically prune redundant tokens. Furthermore, a ranking-based optimization strategy is proposed to minimize the ranking divergence of the saliency score and the predicted importance score. Extensive experiments have shown that our framework is generalizable to various models and datasets. By hierarchically pruning 65% of the input tokens, our method greatly reduces 33% ∼ 47% FLOPs and achieves speedup up to $1.75\times$ during inference, while maintaining comparable performance. We further demonstrate that SDTP can be combined with KV cache compression method for further compression.

## 1 Introduction

Due to the quadratic complexity of the attention mechanism, the efficiency of a large language model is highly influenced by extended input and output sequence length. Using Mistral-7B as an example, a model with 7 billion parameters. When the input sequence length is increased from 4K tokens to 128K tokens, the TFLOPs required for processing surges from 72.51 to 8864.49, reflecting a 122-fold rise. This significant growth highlights the substantial computational demands imposed by longer sequences on large language models. To overcome this limitation, token or context pruning turns out to be an important research aspect. On the one hand, precise pruning of less informational tokens directly reduces the input size and the computation FLOPs involved. On the other hand, different from neural network pruning Liu et al. [2023], Fan et al. [2019], Ma et al. [2023], Frantar and Alistarh [2023], distillation Gu et al. [2023], Jiao et al. [2019], Sanh et al. [2019] or quantization Frantar et al. [2022], Xiao et al. [2023a], Wu et al. [2023], it avoids the modification of the original model deployment, therefore the latency improvement is guaranteed across different hardware deployments Rao et al. [2021].

Despite the conceptual advancement, the current token pruning methods still face challenges in LLM inference. Firstly, most methods focus on KV cache compression during autoregressive decoding Zhang et al. [2024], Liu et al. [2024], Anagnostidis et al. [2024], Xiao et al. [2023b]. While these methods are generally beneficial to LLM inference speed, they are not designed for
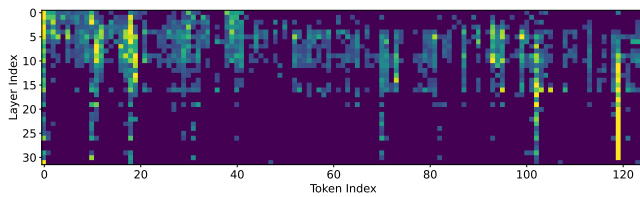
---

\*Corresponding Author.
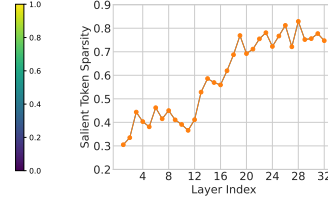
Figure 1: Gradient-based saliency scores.

Figure 2: Salient token sparsity.

the computation-bounded stage of prefill, leading to a sub-optimal acceleration in many real world scenarios, especially when prompt engineering generates a long in-context content injection for the input Jiang et al. [2023b], Pan et al. [2024]. Moreover, the token selection criteria is often limited to local attention statics (e.g., Zhang et al. [2024]) or sigmoid functions (e.g., Anagnostidis et al. [2024]), which may neglect important semantic information in the pruning process.

To overcome these limitations, in this study, we propose a novel token pruning framework, namely Saliency-driven Dynamic Token Pruning (SDTP), to perform a cascade pruning process across the Transformer layers to speed up inference. The motivation of our method originates from gradient-based model interpretability Kindermans et al. [2019]. Following the intuition that tokens with higher gradient-based saliency scores indicate more linguistic contribution, we first conduct a preliminary analysis of observations on the saliency map of LLM models. Figure 1 shows a typical example using LIama2-7B, where we find that not all input tokens have the same contribution. In particular, we find that a) important tokens are sparse, with a sparse rate increasing with the number of layers; and b) the sparsity can largely be passed to deeper layers, i.e., if the token at the front layer is determined as redundant, it will still be redundant in deeper layers. The sparsity pattern can also be statistically verified in Figure 2 using the saliency statistics on the Dolly dataset. The important tokens are defined as the tokens whose saliency score is greater than 10% of the maximum value, and the rest are redundant tokens.

Our framework focuses on determining the most informative tokens given the input prompt. In particular, we leverage a gradient-based method to dynamically mark important tokens across the layers, and finally generate a pruned input before the first token generation. Our main contributions are summarized as follows:

- We develop a token selection scheme by pruning the input length, which benefits the whole process of LLM inference, achieving nearly lossless speedup even with a high pruning ratio.

- We propose a novel learnable pruning module supervised by gradient-based feature attribution and a ranking-based optimization strategy.

- Extensive experiments have shown our framework is generalizable to various models and datasets by hierarchically pruning 65% of the input tokens, achieving up to 47.2% FLOPs reduction, 1.75× speeding up, and 34.26% memory savings for Mistral-7B.

## 2 Related Work

**Structure pruning for transformer.** In the context of transformer models for natural language processing tasks, a significant body of work has concentrated on reducing the number of heads in the multi-head attention (MSA) module. For instance, Michel et al. Michel et al. [2019] observe that eliminating a large percentage of heads in pre-trained BERT models Devlin et al. [2018] has a minimal impact on performance. Beyond the MSA module, pruning has also been applied to the tokens, as demonstrated in Goyal et al. [2020], Kim and Cho [2020], Wang et al. [2021]. LTP Kim et al. [2022] is a threshold-based token pruning method that needs a simple threshold operation and fully automates the search for optimal pruning configurations with a differentiable soft binarized mask. Yun Yun et al. [2023] proposes an approach of integrating token pruning and token combining to improve document classification for the BERT model, which includes fuzzy-based token pruned attention and a token combining module that gradually removes unimportant tokens and reduces time and memory costs.
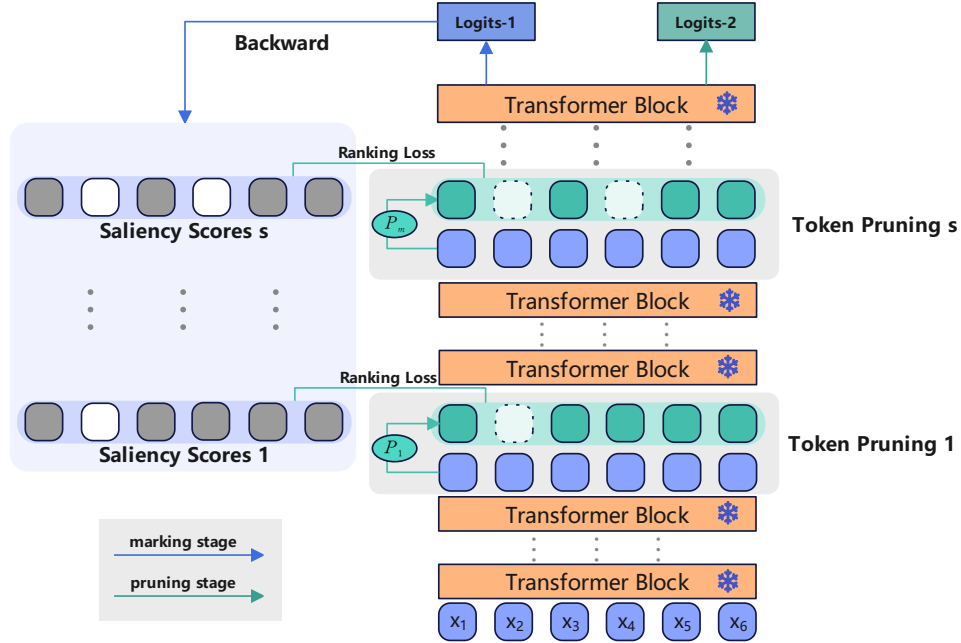
Figure 3: The overall architecture of our proposed method. Our proposed token pruning module is inserted between the transformer blocks and learns to decide the tokens to be pruned. The pruned tokens lead to less computation demand in the following layers.

**Structure pruning for LLMs.** Recent approaches Fan et al. [2019], Kurtic et al. [2023], Liu et al. [2024] employ structural pruning methods on large language models (LLMs), removing entire structured components to achieve more efficient GPU speedups. Based on the contextual sparsity hypothesis, Deja Vu Liu et al. [2023] trims specific attention heads and MLP parameters without modifying pre-trained models. H2O Zhang et al. [2024] is a KV cache eviction policy that identifies a subset of tokens which contribute the most value to attention scores. This method reduces memory footprint and improves LLM deployment efficiency, enhancing throughput and reducing latency without compromising generation quality. LLMLingua Jiang et al. [2023b] proposes a data distillation method to compress prompts from GPT-4 Achiam et al. [2023] by treating prompt compression as a token classification task, capturing essential information, reducing latency, and ensuring the compressed prompt's faithfulness to the original content. Unlike them, our proposed method investigates redundancy from a novel perspective by considering the information integration of different tokens in a transformer. Additionally, reducing tokens can be combined with pruning in other dimensions to achieve a faster acceleration.

## 3 Method

### 3.1 Overview

Our method progressively identifies and removes non-informative tokens from the model's memory during inference. We achieve this by introducing learnable pruning modules within the Transformer architecture. These modules dynamically analyze token representations and select tokens for pruning. Each pruning decision leverages a corresponding gradient-based signal, guided by a dual-loss strategy detailed later. This strategy balances the token pruning objective with the preservation of the model's pre-trained capabilities through a language modeling loss. The pruning mask, indicating the tokens to be removed, is propagated through subsequent layers. Finally, all identified non-informative tokens

are eliminated from memory before calculating the first token. This selective removal significantly reduces computational demands during inference, as the pruning module operates with a fraction of the FLOPs required by the entire language model (less than 1%).

## 3.2 Learnable Token Pruning Module

The architecture of our SDTP is shown in Figure 3. Our SDTP contains a typical transformer and a few token pruning modules. The transformer can be implemented as various large language models such as Llama2 Touvron et al. [2023], BLOOM Le Scao et al. [2023] and Mistral Jiang et al. [2023a]. The token pruning module decides which tokens are dropped by producing the probabilities for all the tokens. Tokens with a smaller score are dropped. We hierarchically insert the token pruning module through the transformer network at certain locations. Thus, the token sparsification is only conducted in these locations. We adopt a two-stage approach in training. The first stage is the token marking stage (blue line), using the original model and all the tokens to perform inference and compute the saliency scores using the output of the transformer. The saliency scores are further utilized in the second stage for computing losses. The second stage is called the token pruning stage (green line). The token pruning module takes effect. We supervise the learning of the token pruning module using MSE loss and a proposed ranking loss. A cross-entropy loss is also utilized to train the transformer. The token pruning module and the transformer can be optimized end-to-end. During inference, the most important tokens are selected and kept according to a predefined ratio in each location, and the scores are computed by the token pruning module.

In the second stage, we use a binary decision mask $M$ to identify the tokens to be kept and pruned. $N$ is the number of tokens. All elements in the decision mask are firstly set to 1 and we update the mask progressively. The token pruning modules take the existing tokens $x \in \mathbb{R}^N$ as input and outputs the decision mask. The token pruning module is implemented as a two-layer MLP:

$$\pi = \text{MLP}(\text{GELU}(\text{MLP}(\mathbf{x}))), \tag{1}$$

$$M = \text{Gumbel-Softmax}(\pi), \tag{2}$$

where $\pi \in \mathbb{R}^{N \times 2}$ is the output of the token pruning module. $M \in \{0, 1\}^N$ is a one-hot tensor output from Gumbel-Softmax. $M$ denotes the mask of tokens that need to be kept. Gumbel-Softmax is a differentiable function, which makes it possible to conduct end-to-end training. GELU is the activation function.

## 3.3 Token Pruning Supervision

Gradient-based attribution techniques Kindermans et al. [2019], Mohebbi et al. [2021] can be used to determine the importance of each input feature through the analysis of partial derivatives of the output with respect to each input dimension. The magnitude of the derivatives serves to reflect the sensitivity of the output with respect to alterations in the input. We use this method to identify the important tokens and further supervise the training of the token pruning module:

$$\hat{\pi} = \frac{\partial T(x)}{\partial x} \cdot x, \tag{3}$$

where $T(x)$ is the output of the entire network and $x$ denotes the input vector of the selected layer that has a token pruning module. $\hat{\pi}$ indicates the importance of tokens in the transformer. Since the saliency score $\hat{\pi}$ cannot be computed in the inference stage, we utilize a token pruning module to simulate it. We use MSE loss to minimize the difference between the predictions of the token pruning module and saliency score $\hat{\pi}$.

$$\mathcal{L}_{\text{mse}} = \text{MSE}\left(\pi, \hat{\pi}\right), \tag{4}$$

The MSE loss ensures the value of $\pi$ and $\hat{\pi}$ become closer. Beyond the value proximity, a more significant consideration is the proximity in the ranking of token importance, since we keep the tokens with higher importance values and a predefined ratio of tokens with smaller importance values are pruned. To minimize the ranking divergence of $\pi$ and $\hat{\pi}$, we propose a ranking-based token selection

loss to supervise the learning of token pruning, as shown in Figure 3. The total ranking loss $\mathcal{L}_r$ is the sum of ranking loss in each stage. For simplicity, we formulate this process as follows:

$$\mathcal{L}_r = \sum_{s=1}^{S} \mathcal{L}_r^{(s)}, \tag{5}$$

where $S$ is the largest number of selected stages that need to be pruned in the transformer. The ranking loss in each layer can be formulated as follows:

$$\mathcal{L}_r^{(s)}(\pi, \hat{\pi}) = \sum_{i=1}^{N-1} \sum_{j=i+1}^{N} log(1 + e^{-((\pi_i - \pi_j) * sign(\hat{\pi}_i - \hat{\pi}_j))}), \tag{6}$$

where $N$ is the total number of tokens. The cross-entropy loss is also utilized to train the entire transformer:

$$\mathcal{L}_{cls} = \text{CrossEntropy}(y, \hat{y}), \tag{7}$$

where $\hat{y}$ is the output of the transformer and $y$ is the label. The total loss is formulated as follows:

$$\mathcal{L} = \mathcal{L}_{cls} + \mathcal{L}_{mse} + \mathcal{L}_r. \tag{8}$$

Our method will be trained end-to-end. We can also distill the knowledge from the original model to the pruned model in training, which improves the performance of the pruned model.

## 4 Experiments

In this section, we aim to demonstrate that SDTP, a lightweight saliency-driven dynamic token pruning framework, can reduce memory usage and latency in wall-clock while maintaining generation quality across a broad spectrum of domains and tasks. Experiments on different LLMs show that the proposed method can be applied to all current mainstream Transformer-based LLMs. Ablation experiments confirm the effectiveness of our proposed SDTP framework and the ingenious ranking-based optimization strategy. We further illustrate the compatibility of our approach with KV cache compression techniques.

### 4.1 Experimental Setup

**Implementation details.** We conducted experiments using three representative LLM models, including the Llama2-7B[2], Mistral-7B-Instruct-v0.2[3] and BLOOM-7B[4]. Unless otherwise specified, we fix the number of pruning stages $S = 10$ and apply the target keeping ratio $r$ as a geometric sequence $[r, r^2, ..., r^{10}]$ where $r = 0.9$.

During the training phase, the pre-trained model is used for initialization, the weight of the original model is frozen, and only the weight of the SDTP module is iteratively updated. This brings two benefits: on the one hand, the strong generalization capability of LLM on downstream tasks can be retained. On the other hand, we only need low-resource alignment training to achieve excellent performance. In addition, for the stability of training, we use attention masking to simulate token pruning, just like in DynamicVit Rao et al. [2021]. All models are trained on the databricks-dolly-15K[5] dataset for two epochs.

During the inference phase, referring to StreamLLM Xiao et al. [2023b] and $H_2O$ Zhang et al. [2024], we always keep the first four initial tokens and the local 10% of the tokens. We set the batch size to 1 to test the accuracy and speed. We use the same hardware for training and evaluation as this work Shao et al. [2024].

---

[2]https://llama.meta.com/

[3]https://mistral.ai/

[4]https://bigscience.huggingface.co/blog/bloom

[5]https://github.com/databrickslabs/dolly/tree/master

**Evaluation of generality.** In order to evaluate the generality of our SDTP on different downstream tasks, we employed the lm-eval-harness Gao et al. [2021] to compare the accuracy before and after dynamic token pruning on several few-shot downstream tasks, including COPA Roemmele et al. [2011], PIQA Bisk et al. [2020], Winogrande Sakaguchi et al. [2021], MathQA Amini et al. [2019], BoolQ Clark et al. [2019], CB Roberts et al. [2020], WiC Pilehvar and Camacho-Collados [2018], and WSC Kocijan et al. [2020]. The evaluation metrics for each task refer to lm-eval-harness Gao et al. [2021].

**Evaluation of long context.** For long-context scenarios, we conduct experiments on LongBench Bai et al. [2023] with Mistral-7B-Instruct-v0.2[6], which achieves compelling performance in long-context scenarios. We evaluate our method and LLMLingua-2 Pan et al. [2024] on the English datasets of LongBench, including single-document QA, multi-document QA, summarization, few-shot learning and code generation tasks. We prune the model with 5 and 10 pruning stages, denoted as SDTP-low and SDTP in Table 1, respectively. The compression ratio is the number of original tokens divided by the number of compressed tokens. In our method, the average number of input tokens of all layers is taken as the number of compressed tokens. The evaluation metrics for each dataset refer to LongBench github[7]. We evaluate the metrics and inference latency with Huggingface's transformers 4.36.0, PyTorch 2.0.1, and FlashAttention2 2.3.3.

## 4.2 Main Results

**Results on long context benchmarks.** Considering that there is no related work to dynamically prune tokens in the middle layers during inference in the LLM scenario, we compared our method with the prompt pruning method LLMLingua-2 Pan et al. [2024] in the out-of-domain scenario, which is a strong baseline. As illustrated in Table 1, with the same token compression ratio of 1.61x, our method SDTP achieves higher performance and speedup ratio on LongBench than LLMLingua-2. The average performance of our method is 4.66 higher than LLMLingua-2, as SDTP performs significantly better on both Code and FewShot tasks. The speedup ratio gain is attributed to our lightweight pruning module, while LLMLingua-2 uses xlm-roberta-large and multilingual-BERT as prompt compressors. These results demonstrate that our method can prune less important tokens more efficiently. To reduce the performance gap between our method and the original input tokens, SDTP-low is pruned with a lower compression ratio of 1.32x, and the performance drop on every task is reduced to 1~2.37.

Table 1: Evaluation with Mistral-7B on LongBench. SDTP-low is pruned with a lower compression ratio than SDTP. $1/\tau$ is the compression ratio. * denotes our implementation based on the official codes.

| Methods | LongBench | | | | | | $1/\tau$ | Speedup | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | SingleDoc | MultiDoc | Summ. | FewShot | Code | AVG | | Prefill | End2End |
| Original Token | 36.66 | 29.88 | 28.26 | 66.95 | 53.74 | 42.34 | - | - | - |
| LLMLingua-2-small* | 31.55 | 29.04 | 26.79 | 49.24 | 32.61 | 33.94 | 1.61× | 1.16× | 1.12× |
| LLMLingua-2* | 31.35 | **30.57** | **27.29** | 52.29 | 33.87 | 35.16 | 1.61× | 0.83× | 0.97× |
| SDTP | **32.70** | 27.74 | 26.65 | **62.74** | **53.99** | **39.82** | 1.61× | **1.71×** | **1.31×** |
| SDTP-low | 35.60 | 28.83 | 27.26 | 64.58 | 52.43 | 40.98 | 1.32× | 1.37× | 1.10× |

**Inference latency and FLOPs.** In order to analyze the performance gains under different sequence lengths in detail, a series of speed tests were conducted using specifically constructed data. In these tests, the prompt of the constructed data is padded to the same length, and the LLMs are required to generate tokens of the same length, which is set to 128 here. The context lengths of the prompts ranged from 4K to 128K. The theoretical maximum reduction in end-to-end FLOPs after the integration of the SDTP module was predicted to be 47.2%. The end-to-end and prefill phase inference latency were recorded for various prompt lengths. Additionally, the peak GPU memory usage during inference was monitored.

---

[6]https://mistral.ai/
[7]https://github.com/THUDM/LongBench

Table 2: Inference latency, FLOPs and GPU memory usage test on Mistral-7B with context lengths from 4K to 128. Our SDTP can achieve up to 47.2% FLOPs reduction, 1.75× end-to-end acceleration, and 34.26% memory savings.

| Seq. length | Method | TFLOPs | | Latency(s) | | Memory(GB) |
|---|---|---|---|---|---|---|
| | | Prefill | E2E | Prefill | End2End | |
| 4K | w/o SDTP | 72.51 | 74.78 | 0.31 | 3.70 | 15.55 |
| | w SDTP | 47.81 | 50.01 | 0.20(↑ 1.52×) | 3.43(↑ 1.08×) | 15.04(↓ 3.28%) |
| 8K | w/o SDTP | 158.22 | 160.70 | 0.64 | 4.85 | 17.10 |
| | w SDTP | 102.01 | 104.35 | 0.40(↑ 1.58×) | 3.98(↑ 1.22×) | 15.86(↓ 7.25%) |
| 16K | w/o SDTP | 369.22 | 372.11 | 1.40 | 7.17 | 20.20 |
| | w SDTP | 229.41 | 232.03 | 0.84(↑ 1.67×) | 5.46(↑ 1.31×) | 17.55(↓ 13.11%) |
| 32K | w/o SDTP | 949.54 | 953.25 | 3.40 | 13.63 | 26.41 |
| | w SDTP | 560.40 | 563.56 | 1.90(↑ 1.79×) | 8.95(↑ 1.52×) | 20.95(↓ 20.67%) |
| 64K | w/o SDTP | 2743.46 | 2748.83 | 9.33 | 27.75 | 39.31 |
| | w SDTP | 1526.98 | 1531.25 | 4.99(↑ 1.87×) | 16.92(↑ 1.64×) | 28.24(↓ 28.18%) |
| 128K | w/o SDTP | 8864.49 | 8873.15 | 29.30 | 63.69 | 65.13 |
| | w SDTP | 4678.81 | 4685.28 | 13.60(↑ 2.15×) | 36.35(↑ 1.75×) | 42.82(↓ 34.26%) |

The results presented in Table 2 indicate that the SDTP module significantly reduces inference latency, achieving an end-to-end speedup of up to 1.75 times. This speedup is substantial, demonstrating the effectiveness of the SDTP module in optimizing model performance. Furthermore, the SDTP module also effectively reduced the GPU memory required during inference, leading to a memory saving of up to 34.26%. This reduction in memory usage is particularly beneficial for large-scale applications with limited computational resources.

Overall, the findings from these experiments confirm that the SDTP module not only enhances the speed of inference but also reduces the memory footprint of LLMs, making them more suitable for deployment on systems with constrained resources.

## 4.3 Generalization of SDTP

**Generalization of SDTP on different tasks.** Since our method freezes all Transformer blocks and only inserts dynamic token pruning modules between blocks, our method can theoretically preserve the generalization of the original LLM model. We have verified this through experiments on three LLMs with eight 5-shot downstream tasks: COPA Roemmele et al. [2011], PIQA Bisk et al. [2020], Winogrande Sakaguchi et al. [2021], MathQA Amini et al. [2019], BoolQ Clark et al. [2019], CB Roberts et al. [2020], WiC Pilehvar and Camacho-Collados [2018], and WSC Kocijan et al. [2020]. As shown in Table 3, the model retains comparable performance to the original model after cumulatively pruning up to 65% of tokens.

**Results of different LLMs.** As shown in Table 3, we performed experiments on three Transformer-based LLMs with 7 billion parameters: Llama2-7B, Mistral-7B, and BLOOM-7B. The results demonstrate that our approach achieves accuracy comparable to all three original, unpruned LLMs. This finding suggests the versatility of the SDTP framework, indicating its potential for integration into various Transformer-based LLMs without sacrificing performance.

## 4.4 Ablation Results

**Ablation on the saliency-driven dynamic token pruning.** For comparison, we implemented a token pruning module with only token pruning rate constraints similar to DynamicVit Rao et al. [2021], and trained it on the Dolly dataset to make it suitable for NLP tasks as our baseline. We test the effectiveness of our proposed saliency-driven token selection module and ranking-based optimization strategy using two target LLM models: Llama2-7B and Mistral-7B. As illustrated in Table 4, we

Table 3: Accuracy of eight 5-shot downstream tasks on three LLMs: Llama2-7B, Mistral-7B, and BLOOM-7B. Our SDTP is generalizable to various models and tasks by hierarchically pruning 65% of the input tokens.

| Method | Llama2-7B | | Mitral-7B | | BLOOM-7B | |
|---|---|---|---|---|---|---|
| | Full | SDTP | Full | SDTP | Full | SDTP |
| COPA | 83.00 | 85.00 | 90.00 | 91.00 | 69.00 | 69.00 |
| PIQA | 78.84 | 78.89 | 84.71 | 84.44 | 73.61 | 71.60 |
| WinoGrande | 73.88 | 72.22 | 76.64 | 76.16 | 65.59 | 60.85 |
| MathQA | 31.93 | 31.19 | 36.45 | 36.78 | 27.81 | 27.37 |
| BoolQ | 60.98 | 66.42 | 66.91 | 71.35 | 61.19 | 61.77 |
| CB | 51.79 | 53.57 | 92.86 | 89.29 | 48.21 | 53.57 |
| WiC | 52.04 | 55.17 | 57.21 | 55.49 | 50.31 | 50.47 |
| WSC | 36.54 | 36.54 | 66.35 | 65.38 | 36.54 | 36.54 |
| Avg | 58.62 | 59.88 | 71.39 | 71.24 | 54.03 | 53.90 |
| Pruning ratio | - | 65% | - | 65% | - | 65% |

Table 4: Ablation on four 5-shot tasks with Llama2-7B as target LLM. Our proposed saliency-driven token selection module and ranking-based optimization strategy can improve the accuracy of downstream tasks, with an improvement of 0.89 on average compared to the baseline.

| Method | COPA | PIQA | Winogrande | MathQA | Average |
|---|---|---|---|---|---|
| SDTP | **85.00** | **78.89** | **72.22** | **31.19** | **66.82** |
| $w/o$ Rank | 84.00(-1.00) | 78.40(-0.49) | 71.82(-0.40) | 30.85(-0.34) | 66.27(-0.55) |
| $w/o$ Rank & MSE | 83.00(-2.00) | 78.29(-0.60) | 71.43(-0.79) | 30.99(-0.20) | 65.93(-0.89) |

compared four 5-shot evaluation tasks: COPA, PIQA, Winogrande, MathQA. The results demonstrate that our proposed saliency-driven dynamic token pruning module effectively retains crucial tokens and that pairwise ranking loss further aids the training of the token pruning module. Our SDTP can improve the accuracy of downstream tasks with an improvement of 0.89 on average compared to the baseline. As shown in Table 5, we compared five types of tasks, such as single-document QA and multi-document QA on the multi-task long context understanding benchmark LongBench Bai et al. [2023]. The results demonstrate that our SDTP can improve the accuracy of long context understanding tasks, especially for multi-document QA tasks, which can be improved by 1.12.

Table 5: Ablation on LongBench with Mistral-7B as target LLM. Our proposed saliency-driven token selection module and ranking-based optimization strategy can improve the accuracy of long context understanding tasks. Especially on multi-document QA tasks, it can be improved by 1.12.

| Method | SingleDoc | MultiDoc | Summ. | FewShot | Code |
|---|---|---|---|---|---|
| SDTP | **32.65** | **27.28** | **26.60** | **62.63** | **53.81** |
| $w/o$ Rank | 32.23(-0.22) | 26.51(-0.77) | 26.32(-0.28) | 61.91(-0.72) | 53.34(-0.47) |
| $w/o$ Rank & MSE | 31.94(-0.71) | 26.16(-1.12) | 26.59(-0.01) | 62.00(-0.63) | 53.59(-0.22) |

## 4.5 Compatible with KV Cache Compression

The token pruning method is designed to reduce the number of FLOPs during the prefill phase, while the KV cache compression method aims to decrease FLOPs in the decode phase. Since these two methods target different stages of the computation and do not interfere with each other, they are considered orthogonal. In the experiment described, the compatibility of the proposed method with KV cache compression was tested with Llama2-7B. The $H_2O$ Zhang et al. [2024] algorithm was employed for this purpose. During the KV cache storage phase, $H_2O$ compresses the KV cache, reserving 40% of it, which can lead to further improvements in inference speed and throughput. Table 6 indicates that the proposed method is compatible with KV cache compression, as it does not introduce compound errors. This means that the benefits of both techniques reduced FLOPs from

Table 6: Compatibility with KV cache compression with Llama2-7B as target LLM. Our SDTP is compatible with KV cache compression and does not introduce compound errors by hierarchically pruning 65% of the input tokens.

| Tasks | Full | 40% KV | | |
| --- | --- | --- | --- | --- |
| | | Local | $H_2O$ | SDTP $w.H_2O$ |
| COPA | 83.00 | 50.00 | 85.00 | 85.00 |
| PIQA | 78.84 | 50.27 | 78.62 | 78.56 |
| Winogrande | 73.88 | 48.46 | 73.32 | 72.53 |
| MathQA | 31.93 | 23.38 | 30.99 | 30.99 |
| Pruning ratio | - | - | - | 65% |

token pruning and efficient cache management from compression can be simultaneously realized without negatively affecting the accuracy or performance of the model.

## 4.6 Experiment on Pruning Rate and Pruning Strategy

To optimize the token pruning strategy for minimal impact on model performance, a series of experiments were conducted using Llama2-7B as the target LLM on the 5-shot task Winogrande. These experiments involved single-stage pruning from the first layer up to the tenth layer, with varying keeping ratios of 90%, 80%, and 70%. The goal was to identify the optimal layer to start pruning and the ideal keeping ratio that would least affect the model's performance. As shown in Figure 4, it was observed that starting the pruning process from the fourth layer with a conservative keeping ratio of 90% can achieve a balance between model efficiency and performance maintenance. This finding suggests that earlier layers are less amenable to aggressive pruning. Conversely, later layers, which tend to be more specialized and may contain more redundant or less critical information, can tolerate higher levels of pruning without significant degradation in performance. This outcome is important for designing efficient pruning strategies that reduce computational overhead without compromising accuracy.

To further investigate the effect of multi-stage pruning on model accuracy, experiments were conducted across eight 5-shot tasks. In these experiments, a keeping ratio of 90% was maintained while pruning began from the fourth layer, with a layer step of 3. The layer step refers to the number of layers between each successive pruning stage. As shown in Figure 5, the results indicated that although a high pruning rate in a single stage could be detrimental to model performance, increasing the number of pruning stages had a minor impact on performance. This is in line with the earlier observation in the article regarding the increase in token sparseness with deeper layers. The rule suggests that significant token sparseness occurs as you move deeper into the layers of LLM, allowing for more aggressive pruning in later stages without substantially affecting the model's ability to process and generate language.
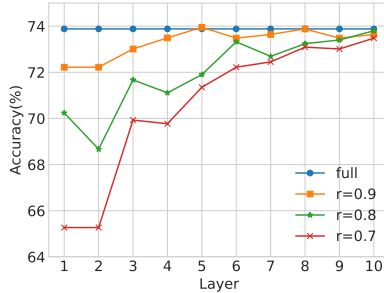


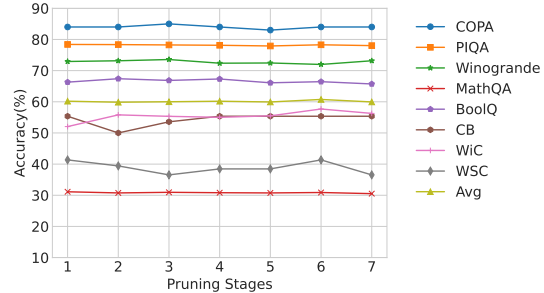Figure 4: Effect of first pruning layer.



Figure 5: Effect of multi-stage pruning.

9

# 5   Conclusion

This work introduces a novel approach to accelerate LLM inference through dynamic redundant token pruning. We propose the Saliency-Driven Token Pruning (SDTP) module, a lightweight and adaptable component that seamlessly integrates within the Transformer layer. The SDTP module leverages a supervised training strategy guided by token saliency and a ranking-based optimization strategy to effectively select tokens for pruning during inference. Our experimental results demonstrate that SDTP achieves significant inference acceleration, progressively pruning up to 65% of input tokens while attaining a speedup of $1.75\times$. Notably, this method is applicable to a wide range of mainstream Transformer-based LLMs and offers the potential for further compression when combined with KV cache compression techniques.

# References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

Aida Amini, Saadia Gabriel, Peter Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. Mathqa: Towards interpretable math word problem solving with operation-based formalisms. *arXiv preprint arXiv:1905.13319*, 2019.

Sotiris Anagnostidis, Dario Pavllo, Luca Biggio, Lorenzo Noci, Aurelien Lucchi, and Thomas Hofmann. Dynamic context pruning for efficient and interpretable autoregressive transformers. *Advances in Neural Information Processing Systems*, 36, 2024.

Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, et al. Longbench: A bilingual, multitask benchmark for long context understanding. *arXiv preprint arXiv:2308.14508*, 2023.

Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439, 2020.

Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint arXiv:1905.10044*, 2019.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

Angela Fan, Edouard Grave, and Armand Joulin. Reducing transformer depth on demand with structured dropout. *arXiv preprint arXiv:1909.11556*, 2019.

Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot. In *International Conference on Machine Learning*, pages 10323–10337. PMLR, 2023.

Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022.

Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, et al. A framework for few-shot language model evaluation. *Version v0. 0.1. Sept*, page 8, 2021.

Saurabh Goyal, Anamitra Roy Choudhury, Saurabh Raje, Venkatesan Chakaravarthy, Yogish Sabharwal, and Ashish Verma. Power-bert: Accelerating bert inference via progressive word-vector elimination. In *International Conference on Machine Learning*, pages 3690–3699. PMLR, 2020.

Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. Minillm: Knowledge distillation of large language models. In *The Twelfth International Conference on Learning Representations*, 2023.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023a.

Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. Llmlingua: Compressing prompts for accelerated inference of large language models. *arXiv preprint arXiv:2310.05736*, 2023b.

Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. Tinybert: Distilling bert for natural language understanding. *arXiv preprint arXiv:1909.10351*, 2019.

Gyuwan Kim and Kyunghyun Cho. Length-adaptive transformer: Train once with length drop, use anytime with search. *arXiv preprint arXiv:2010.07003*, 2020.

Sehoon Kim, Sheng Shen, David Thorsley, Amir Gholami, Woosuk Kwon, Joseph Hassoun, and Kurt Keutzer. Learned token pruning for transformers. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 784–794, 2022.

Pieter-Jan Kindermans, Sara Hooker, Julius Adebayo, Maximilian Alber, Kristof T Schütt, Sven Dähne, Dumitru Erhan, and Been Kim. The (un) reliability of saliency methods. *Explainable AI: Interpreting, explaining and visualizing deep learning*, pages 267–280, 2019.

Vid Kocijan, Thomas Lukasiewicz, Ernest Davis, Gary Marcus, and Leora Morgenstern. A review of winograd schema challenge datasets and approaches. *arXiv preprint arXiv:2004.13831*, 2020.

Eldar Kurtic, Elias Frantar, and Dan Alistarh. Ziplm: Hardware-aware structured pruning of language models. *arXiv preprint arXiv:2302.04089*, 3(7), 2023.

Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-parameter open-access multilingual language model. 2023.

Zichang Liu, Jue Wang, Tri Dao, Tianyi Zhou, Binhang Yuan, Zhao Song, Anshumali Shrivastava, Ce Zhang, Yuandong Tian, Christopher Re, et al. Deja vu: Contextual sparsity for efficient llms at inference time. In *International Conference on Machine Learning*, pages 22137–22176. PMLR, 2023.

Zichang Liu, Aditya Desai, Fangshuo Liao, Weitao Wang, Victor Xie, Zhaozhuo Xu, Anastasios Kyrillidis, and Anshumali Shrivastava. Scissorhands: Exploiting the persistence of importance hypothesis for llm kv cache compression at test time. *Advances in Neural Information Processing Systems*, 36, 2024.

Xinyin Ma, Gongfan Fang, and Xinchao Wang. Llm-pruner: On the structural pruning of large language models. *Advances in neural information processing systems*, 36:21702–21720, 2023.

Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? *Advances in neural information processing systems*, 32, 2019.

Hosein Mohebbi, Ali Modarressi, and Mohammad Taher Pilehvar. Exploring the role of bert token representations to explain sentence probing results. *arXiv preprint arXiv:2104.01477*, 2021.

Zhuoshi Pan, Qianhui Wu, Huiqiang Jiang, Menglin Xia, Xufang Luo, Jue Zhang, Qingwei Lin, Victor Rühle, Yuqing Yang, Chin-Yew Lin, et al. Llmlingua-2: Data distillation for efficient and faithful task-agnostic prompt compression. *arXiv preprint arXiv:2403.12968*, 2024.

Mohammad Taher Pilehvar and Jose Camacho-Collados. Wic: the word-in-context dataset for evaluating context-sensitive meaning representations. *arXiv preprint arXiv:1808.09121*, 2018.

Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. Dynamicvit: Efficient vision transformers with dynamic token sparsification. *Advances in neural information processing systems*, 34:13937–13949, 2021.

Adam Roberts, Colin Raffel, and Noam Shazeer. How much knowledge can you pack into the parameters of a language model? *arXiv preprint arXiv:2002.08910*, 2020.

Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In *2011 AAAI Spring Symposium Series*, 2011.

Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.

Victor Sanh, L Debut, J Chaumond, and T Wolf. Distilbert, a distilled version of bert: Smaller, faster, cheaper and lighter. arxiv 2019. *arXiv preprint arXiv:1910.01108*, 2019.

Ninglu Shao, Shitao Xiao, Zheng Liu, and Peitian Zhang. Flexibly scaling large language models contexts through extensible tokenization. *arXiv preprint arXiv:2401.07793*, 2024.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

Hanrui Wang, Zhekai Zhang, and Song Han. Spatten: Efficient sparse attention architecture with cascade token and head pruning. In *2021 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, pages 97–110. IEEE, 2021.

Xiaoxia Wu, Cheng Li, Reza Yazdani Aminabadi, Zhewei Yao, and Yuxiong He. Understanding int4 quantization for transformer models: Latency speedup, composability, and failure cases. *arXiv preprint arXiv:2301.12017*, 2023.

Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning*, pages 38087–38099. PMLR, 2023a.

Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming language models with attention sinks. *arXiv preprint arXiv:2309.17453*, 2023b.

Jungmin Yun, Mihyeon Kim, and Youngbin Kim. Focus on the core: Efficient attention via pruned token compression for document classification. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023.

Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yuandong Tian, Christopher Ré, Clark Barrett, et al. H2o: Heavy-hitter oracle for efficient generative inference of large language models. *Advances in Neural Information Processing Systems*, 36, 2024.