

FiRST: Finetuning Router-Selective Transformers for Input-Adaptive Latency Reduction

Akriti Jain*, Saransh Sharma*[†], Koyel Mukherjee, Soumyabrata Pal

Adobe Research, India

{akritij, sarsharma, komukher, soumyabratap}@adobe.com

Abstract

Auto-regressive Large Language Models (LLMs) demonstrate remarkable performance across different domains such as vision and language tasks. However, due to sequential processing through multiple transformer layers, autoregressive decoding faces significant computational challenges, particularly in resource-constrained environments like mobile and edge devices. Existing approaches in literature that aim to improve latency via skipping layers have two distinct flavors: (1) early exit, and (2) input-agnostic heuristics where tokens exit at pre-determined layers irrespective of input sequence. Both the above strategies have limitations, the former cannot be applied in the presence of KV caching, which is essential for speed-ups in modern inference frameworks, and the latter fails to capture variation in layer importance across tasks or, more generally, across input sequences. To address these limitations, we propose FiRST, a model-agnostic framework that reduces inference latency by using layer-specific routers to adaptively skip transformer layers during decoding, based on routing decisions made from the input prompt in the prefill stage. FiRST remains fully compatible with KV caching, enabling faster decoding while maintaining quality. Our method reveals that input adaptivity is essential: Different tasks rely on different subsets of layers to evolve meaningful representations. Extensive experiments show that FiRST significantly reduces latency while outperforming existing layer selection strategies in quality. It retains performance comparable to the base model without skipping. FiRST is thus a promising and efficient solution for LLM deployment in low-resource environments.

1 Introduction

Large Language Models (LLMs) have revolutionized the fields of Natural Language Processing and

Computer Vision achieving incredible performance on a diverse set of benchmark tasks. However, the massive scale of LLMs, often involving billions of parameters, poses significant challenges for deployment in resource-constrained environments, where memory, compute, and especially latency become critical bottlenecks. In this work, we focus on addressing the latency issue, which is particularly pronounced in edge settings such as laptops and mobile devices. As noted by Schuster et al. (2022a), the auto-regressive nature of decoding in LLMs further amplifies this bottleneck.

Transformer-based LLMs consist of multiple stacked layers, including attention and feed-forward networks, which result in high latency and computational cost. This makes inference slow or even impractical in resource-constrained environments. The inefficiency stems from the need to process each input token sequentially through all layers, irrespective of the input sequence or task. However, it is important to note that in the real world, there is a lot of heterogeneity in input sequences and tasks. (Schuster et al., 2022a; Sun et al., 2022) observed that LLM generations can have varying levels of difficulty, and certain generations can be solved with reduced computation by exiting the transformer stack early. At the same time, it has been noted in recent works (Wendler et al., 2024) that inference forward pass proceeds in phases through the layers of transformer-based models, with different types of information being extracted or mapped at different phases (sequences of layers) for certain tasks such as translation. Motivated by these and other related works, we hypothesize that *different sequential combinations of layers are important for different input sequences and tasks*. Learning the right sequential combination of layers can help reduce inference latency and compute for on-device scenarios. However, there are several challenges. Any algorithm for determining the “right” combination of layers should minimize

*Equal Contribution

[†]Work done during internship at Adobe Research

any quality loss, be compatible with other latency reduction strategies such as KV cache handling, and be learnable with minimal compute/training overhead.

In the last few years, several promising approaches have been proposed in literature that adaptively prune layers at each decoding step. Token-level early exit proposed in (Schuster et al., 2022a; Sun et al., 2022) allow tokens to exit the transformer layer stack early based on different strategies to compute the confidence or saturation level. (Elhoushi et al., 2024; Elbayad et al., 2020; Zhang et al., 2019) extended this idea to incorporate layer skipping at a token level during training. While token level early exit is a useful idea in theory, it suffers from a major limitation of incompatible KV caching in practice (Del Corro et al., 2023). The incompatibility stems from having to recompute KV caches for preceding tokens if we have a delayed exit point for latter tokens, often resulting in loss of early exit advantages. This limits its practical adoption since KV cache is crucial in significantly speeding up auto-regressive decoding.

Recently, (Liu et al., 2024; Del Corro et al., 2023; Song et al., 2024) have proposed input-agnostic layer skipping at token level, that handle KV cache appropriately as well as retain the advantage of adaptive partial computation. In these solutions, tokens exit at pre-determined layers irrespective of the input sequence, and for all sequences in a batch, tokens at the same position in a sequence exit at the same layer. Furthermore, tokens at latter parts of the sequence are constrained to exit earlier than the previous tokens to ensure that there is no redundant KV cache re-computation. These solutions are heuristic based and impose hard rules and constraints irrespective of input sequences, which can lead to drop in output quality. Others (Jaiswal et al., 2024; Chen et al., 2024) have proposed skipping layers by identifying redundant ones through computing cosine similarity of (input/output) representations of a layer. However, their strategy does not take into account that several middle layers are crucial (see (Liu et al., 2024)) and furthermore, final prediction capability of full model is not taken into account while deciding which layers to skip. *Importantly, in none of the works described above, the strategy of selecting layers for skipping is sequence dependent. Furthermore, they do not consider fine-tuning the models in a way such that not only the performance improves but the model also*

learns to skip layers appropriately.

Due to space constraints, we delegate a study of other related works and orthogonal approaches (for e.g. model compression) for exploring latency/performance tradeoff to Appendix A.1.

Our goal is to design an input-adaptive, **learnable layer selection strategy** that provides quality-aware latency improvements while properly handling KV caching. Ideally, for each input sequence and task, we want to predict a sequential combination of layers to run during inference, minimizing quality loss while achieving as much latency gain as possible. We also want to do this with minimal computational overhead or extra training. To achieve this, we propose training **routers**. At each layer, the router looks at the current sequence representation and decides whether to skip the next layer. Since the decision is made at the sequence level, all tokens follow the same path through the model, avoiding KV cache inconsistencies during decoding. Finally, we fine-tune the model with trained routers using LoRA adapters to recover any quality drop introduced by layer skipping, while preserving latency gains. LoRA fine-tuning also smoothens layer skipping and further highlights the varied importance of layers based on input sequence. Our key contributions include:

1. We propose **FIRST**, a training and inference algorithm that uses layer-specific routers to perform input-adaptive layer selection. All tokens in a sequence follow the same selected layers, ensuring compatibility with KV caching and avoiding additional compute or latency overhead. **FIRST** is model-agnostic and can be applied on top of any pre-trained LLM.
2. We introduce a LoRA-based fine-tuning approach on top of router-based layer selection to recover quality while maintaining latency gains. This also encourages smoother and a more stable layer selection.
3. Finally, we conduct extensive experiments with **FIRST** across multiple datasets spanning three distinct tasks: Machine Translation, Summarization, and Question Answering. We evaluate on two open-source model architectures: LLaMA-3-8B and LLaMA-3.2-3B, and show that, for the same target speed-up, **FIRST** significantly improves performance across tasks as compared to baselines.

2 Problem Statement

Our goal is to exploit the heterogeneity in inputs and tasks to selectively use LLM layers in a quality-aware manner for reducing inference latency and compute for on-device constraints. Ideally, we want to select an *optimal* sub-sequence of layers within a transformer architecture for a given input and task, such that the overall latency, as well as expended computation, are both low, while quality is comparable to the un-modified case where every input sequence passes through every layer. For ease of explanation, without loss of generality, we assume the task is same and simply consider an input sequence for describing the problem.

Let us consider an input sequence $\mathcal{X} = \{x_1, x_2, \dots, x_n\}$ with n tokens. Let there be m transformer layers in the model, where the i^{th} transformer layer is represented as the function $\phi_i(\cdot)$. As stated lucidly in (Wendler et al., 2024), \mathcal{X} is first converted to an initial latent representation $\mathcal{H}_0 = \{H_0^1, H_0^2, \dots, H_0^n\}$, where $H_j^0 \in \mathbb{R}^D, \forall j \in [n]$ is a look-up from a learned embedding dictionary corresponding to the j^{th} token. Thereafter, every transformer layer $\phi_i(\cdot)$ operates on the latent vectors \mathcal{H}_i to generate the embedding for the i^{th} layer as follows. For the j^{th} token, the embedding at layer i is computed as follows:

$$H_i^j = H_{i-1}^j + \phi_i(H_{i-1}^1, H_{i-1}^2, \dots, H_{i-1}^j) \quad (1)$$

Let the (gold) output or generated sequence for an input sequence \mathcal{X} that passed through all m layers of the model with full computation be $\mathcal{Y}_{\mathcal{X}}^*$. Our hypothesis is that for a given input sequence (and task), there exists an optimal subsequence of functions $\mathcal{F}_{OPT}(\mathcal{X})$ out of the full sequence $\{\phi_i, i \in [m]\}$ such that the output generated by passing through this subsequence: $\mathcal{Y}_{OPT, \mathcal{X}} \approx \mathcal{Y}_{\mathcal{X}}^*$. More formally, if Q is a quantitative quality measure on \mathcal{Y} , and $\epsilon \rightarrow 0$ is tolerance in deviation in quality from the gold output, then we hypothesize that there exists an optimal subsequence, using the minimum number of layers, $\mathcal{F}_{OPT}(\mathcal{X})$, such that:

$$Q(\mathcal{Y}_{OPT, \mathcal{X}}) \geq (1 - \epsilon)Q(\mathcal{Y}_{\mathcal{X}}^*), \forall \mathcal{X}. \quad (2)$$

The optimality above is with respect to the minimum subsequence of layers that can help achieve the above, to minimize latency while keeping quality unaffected. Note that, the optimal subsequence $\mathcal{F}_{OPT}(\mathcal{X})$ need to obey the same autoregressive computation on previous tokens as given in Equation 1. Hence, any algorithm that determines the

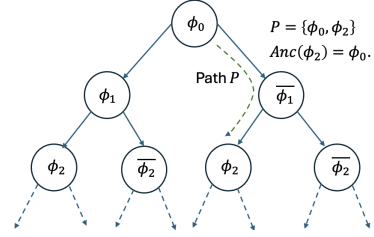


Figure 1: Binary Tree representation of layer selection.

optimal subsequence, need to be compatible with KV cache handling, to avoid the re-computation of values for tokens preceding the current token.

The potential number of subsequences for m layers is 2^m , hence a brute force is infeasible and also beats the purpose of such a layer selection in the first place: reducing latency and compute. In the absence of any known substructure in the behaviour of the latent layers on each input sequence, it is difficult to arrive at the optimal solution polynomially or with low additional latency or compute.

We propose to learn an approximation of the optimal subsequence of layers for any input sequence with low additional latency and minimal training.

3 Proposed Solution: FIRST

Let us first understand what it entails to learn an optimal subsequence of layers for any input. Consider the full transformer sequence to be $\mathcal{F}^* = \{\phi_1, \phi_2, \dots, \phi_m\}$. Any optimal subsequence for an input \mathcal{X} : $\mathcal{F}_{OPT, \mathcal{X}}$ could be thought of as finding an optimal path through a binary tree of functions. Formally, let every level in the binary tree correspond to a transformer layer and the 0^{th} layer corresponds to the initial embedding look up; i.e., at depth $i \in [m]$, there would be 2^i nodes, each corresponding to either ϕ_i or $\overline{\phi}_i$, where the former denotes that a particular transformer layer is included in the optimal path whereas the latter denotes that it is not included. Each (of the 2^{i-1} nodes) ϕ_i or $\overline{\phi}_i$ has two children, corresponding to the next transformer layer: ϕ_{i+1} and $\overline{\phi}_{i+1}$ (See Figure 1). In such a tree structure, for example, the path $\{\phi_i, \overline{\phi}_{i+1}, \phi_{i+2}\}$ indicates the subsequence of transformer layers $\{\phi_i, \phi_{i+2}\}$. For any transformer layer ϕ_i in this tree, let $Anc(\phi_i) = k, 0 \leq k < i$ denote the the lowest ancestor node where the corresponding transformer node ϕ_k is included in the sequence. In the above example, $Anc(\phi_{i+2}) = \phi_i$.

Consider a sequence of functions \mathcal{F} , where for level i , $Anc(\phi_i) = \phi_k$. The autoregressive computations for the j^{th} token in the input sequence

(originally Eq 1), would now be modified as:

$$H_i^j = \begin{cases} H_k^j, & \text{if } \phi_i \notin \mathcal{F}, \\ H_k^j + \phi_i(H_k^1, H_k^2, \dots, H_k^j), & \text{if } \phi_i \in \mathcal{F}. \end{cases} \quad (3)$$

Our problem translates to navigating this binary tree to find the optimal path \mathcal{F}_{OPT} for an input sequence and task. Since there are 2^m paths in this tree, we propose to approximate the optimal by making a decision in a greedy fashion at each node. Formally, we add a (lightweight and fast) router R_i before every transformer layer ϕ_i in the model, that will predict whether ϕ_i will be selected or not.

Our aim is to learn to predict the layer choice at a sequence level (not token) to maintain compatibility with the autoregressive computations and avoid re-computation of of KV cache values. Moreover, we should spend minimal compute for learning the R_i functions. Finally, R_i functions should not add any significant latency to the overall computation.

FIRST modifies any off-shelf pre-trained transformer based model by incorporating and training a router or probability function R_i before every transformer layer ϕ_i . The output of R_i is a score ρ_i denoting the probability of selecting ϕ_i in the layer sequence. During inference, ρ_i is rounded to determine selection of ϕ_i . Let $\lfloor \rho_i \rfloor = 1$ if $\rho_i \geq 0.5$, else 0. Equation 1 is now modified as:

$$H_i^j = H_{i-1}^j + \lfloor \rho_i \rfloor \cdot \phi_i(H_{i-1}^1, H_{i-1}^2, \dots, H_{i-1}^j)$$

This recursively approximates Eq 3 for the optimal \mathcal{F} in a probabilistic, greedy manner. We train the functions R_i on datasets and tasks, and further fine tune using LoRA adapters to make the layer selections smooth and improve the output quality.

4 FIRST Framework and Algorithm

In this section, we describe the training and inference frameworks for FIRST in details. We discuss how to train routers to be adaptive to input sequences. Given an off-the-shelf pre-trained LLM, we propose two training phases. In the first phase, we train a router for each layer that decides whether the input sequence should skip the layer. In the second phase, to tackle the issue of unseen skipping during pre-training, we fine-tune the router-augmented LLM keeping router weights fixed to ensure the model improves performance on the target dataset without reducing the skipping level. While joint training of the router and LoRA modules is

theoretically possible, we find it introduces optimization instability (see Appendix A.4 for further discussion on this design choice and its impact).

4.1 Adaptive Router Module

The adaptive router module is a single-layer neural network without bias, positioned before every layer in the model. During training of the router, all model parameters except the router weights remain frozen. For the first layer, it takes the tokenized input, and for each of the subsequent layers, it takes the output of the preceding layer as input. Mathematically speaking, for any layer i , given a batch of B tokenized inputs sequences, where each sequence has n tokens and is embedded in to \mathbb{R}^D , the adaptive router module takes as input a $B \times n \times D$ tensor output of layer $(i-1)$ and outputs a $B \times n \times 1$ tensor. Subsequently, corresponding to each value (or, token) in the $B \times n \times 1$ tensor, we apply a sigmoid function to ensure that all entries in the tensor are in the interval $[0, 1]$. Following this, we take a mean operation at the sequence level - we take a mean of all the weights in a sequence to output a $B \times 1 \times 1$ tensor. For each sequence in the batch, the corresponding entry is the probability ρ_i with which the sequence passes through the layer i . The input sequence skips the layer i with probability $1 - \rho_i$. During training, the output of a layer is modified using a skip connection, incorporating the probability ρ_i (see Figure: 2).

The routers are trained to encourage skipping by reducing the probabilities $\{\rho_i\}_i$ using a regularizer, to approximate the optimal subsequence for minimizing the latency. The training task is modeled as a language modeling task, specifically next token prediction. The loss function comprises of 3 terms:

- **Cross-entropy loss:** Standard difference between actual and predicted probability distributions to ensure the quality of generation: $\mathcal{L}_{CE} = -\sum_{x \in \mathcal{X}} \mathcal{Y}_x^* \log(\hat{\mathcal{Y}})$.
- **Regularization loss:** Adds a penalty term to reduce overfitting to noise: $\mathcal{L}_{Reg} = \sum_{i \in [m]} \|R_i\|^2$, where $\|R_i\|^2$ denotes the ℓ_2 norm of the router weights for the i^{th} layer router, and there are m layers in the model.
- **Non-skip penalization loss:** This is the summation of probability values across all layers of the model architecture: $\mathcal{L}_{PP} = \sum_{i \in [m]} \rho_i$

The total loss \mathcal{L} is a linear combination of these three terms: $\mathcal{L} = \mathcal{L}_{CE} + \lambda \cdot \mathcal{L}_{Reg} + \alpha \cdot \mathcal{L}_{PP}$, where α manages the tradeoff between quality and latency.

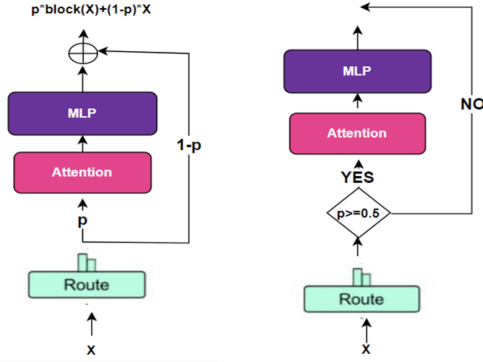


Figure 2: Skip connection used for router training. With probability p , the sequence is processed by the layer and with probability $1 - p$, the layer is skipped. During inference, routers make the decision of whether a sequence will skip a particular layer or pass through it.

4.2 LoRA Compensation Module

Skipping layers naturally leads to some performance loss - especially so since the pre-trained model was not trained to skip layers. To compensate for the loss in performance caused by skipping layers, we finetune the router-augmented pre-trained model on the downstream task¹ using Low Rank Adapters (LoRA). During finetuning, the router parameters are frozen while trainable LoRA adapters are added to both the FFN (Feed-Forward Network) and the attention modules of each layer of the pre-trained model. In order to maintain the skipping level, we again add a non-skip penalization loss component during finetuning with scaling hyper-parameter β . This is essential even though the router weights are frozen because standard finetuning alters the hidden representations of the input sequence in a manner such that no layers are skipped. Note that the LoRA adapters do not lead to any latency overhead during inference.

4.3 Inference for F1RST

During inference, for the input sequence, each router (corresponding to a layer) outputs a number in the interval $[0, 1]$. If this number is greater than or equal to 0.5, the sequence passes through the layer. Otherwise, the sequence skips the layer (Fig. 2). Below, we discuss some salient points about the functioning of the router during inference to handle KV Cache appropriately:

1. **Prefill phase handling:** Skipping is not allowed during prefill phase. This ensures the first token is generated correctly, which is crucial for WMT

¹similar to Quantization Aware Training such as QLoRA (Dettmers et al., 2024) - compensates for model compression

tasks, as they are highly sensitive to the correct generation of the first token in the target language. It has been observed in prior works (Liu et al., 2024) that skipping during prefill phase is detrimental to performance during inference.

2. **Fixed router decisions during decoding and handling KV Cache:** During the prefill phase, the decisions made by the routers are cached. During the decoding phase, every token adheres to the cached decision made during prefill. In other words, for a particular layer, if a router outputs a number less than 0.5 during prefill, the number is fixed for the decoding steps and therefore the same layer will be skipped by all tokens during decoding. Similarly, if the router outputs a number more than 0.5 during prefill, the same layer will be processing all tokens during decoding. Such a step ensures that for each decoding step and each layer that is not skipped, the KV cache for all previous tokens is available for that layer. This approach effectively addresses the caching issues encountered in early exit strategies, ensuring consistent decisions across the decoding process.

5 Experiments

We conduct experiments on three benchmark tasks: Machine Translation, Text Summarization, and Question Answering, demonstrating the robustness and scalability of F1RST across diverse settings. We base our task selection on prior work in the field to ensure a fair and meaningful comparison (Liu et al., 2024; Del Corro et al., 2023; Schuster et al., 2022b).

Datasets: For machine translation, we use WMT development sets (2017–2020) for English-to-Chinese and English-to-German tasks, evaluating performance on the WMT 2022 test set, which covers diverse domains such as news, social media, e-commerce, and conversational contexts. For summarization, we use the CNN/DailyMail dataset, with 4,000 randomly selected training samples and evaluation on the standard test set of 11,490 samples. For question answering (QA), we utilize SQuAD v1.1 and Natural Questions (NQ). We train on 4,000 randomly selected samples from each dataset and evaluate on their respective validation sets as test set labels are unavailable. For NQ, we incorporate a retrieval step before answer generation, retrieving relevant passages as context. Appendix A.2 contains detailed descriptions.

Skip (%)	Model Type		En-to-De				En-to-Zh			
			LLaMA-3-8B		LLaMA-3.2-3B		LLaMA-3-8B		LLaMA-3.2-3B	
			BLEU	COMET	BLEU	COMET	BLEU	COMET	BLEU	COMET
0	Original Model	Base + LoRA	0.199	93.00	0.160	89.72	0.333	82.66	0.278	79.13
		Base	0.169	87.13	0.125	81.66	0.208	68.95	0.166	61.84
15	Skip Decode	Router + LoRA	0.094	55.62	0.105	44.58	0.149	55.98	0.2	46.7
		Router	0.019	23.33	0.055	32.74	0.03	21.75	0.062	34.14
	Random Skip	Router + LoRA	0.097	66.25	0.079	47.26	0.237	67.32	0.154	57.79
		Router	0.132	60.27	0.048	36.3	0.168	59.89	0.077	35.86
	Unified Skip	Router + LoRA	0.153	59.34	0.095	44.72	0.23	69.58	0.157	57.1
		Router	0.117	59.26	0.067	39.81	0.122	54.57	0.087	45.16
	FiRST (Ours)	Router + LoRA	0.161	82.14	0.113	60.29	0.247	68.63	0.218	67.45
		Router	0.108	67.74	0.069	43.04	0.08	42.76	0.1	54.55
25	Skip Decode	Router + LoRA	0.07	31.47	0.077	32.33	0.088	33.85	0.147	42.01
		Router	0.015	21.55	0.05	27.64	0.024	20.93	0.045	29.51
	Random Skip	Router + LoRA	0.018	29.71	0.023	30.97	0.065	27.73	0.137	45.53
		Router	0.011	29.95	0.015	27.22	0.04	35.16	0.053	31.5
	Unified Skip	Router + LoRA	0.06	31.69	0.05	39.81	0.142	50.59	0.108	42.3
		Router	0.048	32.15	0.032	30.86	0.068	38.74	0.079	31.63
	FiRST (Ours)	Router + LoRA	0.071	34.95	0.072	45.38	0.119	56.92	0.126	41.66
		Router	0.029	26.01	0.029	29.08	0.051	25.45	0.053	27.83

Table 1: Machine Translation results for the English-to-German and English-to-Chinese task for both models: BLEU and COMET scores are reported across varying skip levels. **FiRST (Ours)** performs consistently well across both translation directions.

Evaluation Metrics: We employ standard metrics to assess output quality across three distinct tasks. For Machine Translation, we benchmark performance using BLEU and COMET; COMET is included for its more nuanced assessment capabilities beyond the n-gram overlap measured by BLEU. Summarization quality is evaluated with ROUGE scores and BERT Score, the latter capturing meaning-based similarity. For Question Answering, we use Exact Match (EM) and F1 Score on the SQuAD dataset, and report Exact Match (EM) for Natural Questions (NQ). Finally, to benchmark latency, we measure Time Per Output Token (TPOT) on GPU, which gauges overall decoding performance. Detailed descriptions of all evaluation metrics are available in Appendix A.3. Hyperparameters used during training and inference are provided in Appendix A.4.

5.1 Baselines for comparison

We report the latency improvement and quality numbers relative to the base models (no skipping).

- **Random Skipping:** We skip a set of k layers randomly where k depends on the target speedup.
- **Skip Decode:** We implement Skip Decode (Del Corro et al., 2023) method that features a monotonic decrease in processing layers, enabling later tokens to leverage the computational resources used for earlier ones.
- **Unified Skipping:** This, to the best of our knowledge, is the state-of-the-art method relies on using a heuristic-based strategy for retaining layers

at fixed intervals. We replicate the algorithm in (Liu et al., 2024) and compare performance both with and without LoRA fine-tuning across various skipping percentages.

5.2 Experimental Results on Different Tasks

WMT: For LLaMA-3-8B, at 15% skipping, FiRST achieves a latency improvement of upto 10% on TPOT (see Tables 1 and 4). In most cases, it **significantly outperforms** other layer skipping strategies (Skip Decode, Random and Unified Skipping) and in other cases, it is comparable in quality. When compared to the gold output (Base + LoRA), FiRST generally retains a high percentage of the quality, for instance, often achieving $\geq 80\%$ of the COMET score for LLaMA-3-8B, and $\geq 70\%$ in BLEU scores. For 25% skipping, FiRST achieves significant improvement in quality over other strategies, in almost all metrics, while achieving $\sim 18\%$ reduction in TPOT.

For LLaMA-3.2-3B, for similar latency improvement ($\sim 10\%$), the COMET scores are significantly higher than other baselines while the BLEU scores are comparable. (Table 4). It remains within 65 – 85% of BLEU and COMET scores (Table 1) achieved by the gold standard.

CNN/DailyMail Dataset: For LLaMA-3-8B, at $\sim 15\%$ skipping, our method **outperforms** the Base + LoRA setting (Table: 3) while obtaining a 12% improvement in TPOT (Table 4). For LLaMA-3.2-3B, at 15%, the quality is comparable ($\sim 98\%$) to gold and other baselines with 12% improvement in

Skip (%)	Model Type	SQuAD		NQ	
		EM	F1	EM	
LLaMA-3-8B					
0	Original Model	Base + LoRA	73.93	85.99	51.88
		Base	19.46	36.73	37.40
10	Skip Decode	R + LoRA	60.14	65.33	41.98
		Router	16.38	31.48	33.22
	Random Skip	R + LoRA	65.73	80.08	44.95
		Router	18.25	33.75	34.09
	Unified Skip	R + LoRA	55.54	74.58	45.91
		Router	17.39	32.91	33.64
FiRST (Ours)	R + LoRA	70.85	83.61	47.85	
Router		14.58	31.52	33.53	
20	Skip Decode	R + LoRA	45.00	55.10	26.62
		Router	10.68	26.69	14.39
	Random Skip	R + LoRA	47.79	66.37	28.96
		Router	6.71	22.46	27.40
	Unified Skip	R + LoRA	52.87	69.28	25.30
		Router	18.18	32.51	23.57
FiRST (Ours)	R + LoRA	60.60	75.49	32.22	
Router		13.21	27.48	18.02	
LLaMA-3.2-3B					
0	Original Model	Base + LoRA	73.07	84.17	40.50
		Base	18.92	37.74	30.10
10	Skip Decode	R + LoRA	60.79	75.00	31.74
		Router	20.00	31.55	21.51
	Random Skip	R + LoRA	64.78	77.27	36.02
		Router	13.76	28.59	22.76
	Unified Skip	R + LoRA	65.03	77.53	32.90
		Router	13.16	32.31	21.03
FiRST (Ours)	R + LoRA	69.44	81.35	37.82	
Router		12.79	28.37	22.55	
20	Skip Decode	R + LoRA	40.12	40.00	27.28
		Router	20.45	37.62	14.73
	Random Skip	R + LoRA	11.32	38.34	26.60
		Router	6.75	15.51	13.24
	Unified Skip	R + LoRA	37.39	52.49	26.84
		Router	7.81	18.20	16.09
FiRST (Ours)	R + LoRA	39.70	54.59	29.87	
Router		5.52	15.33	17.51	

Table 2: Quality Analysis on Question Answering tasks using LLaMA-3-8B and LLaMA-3.2-3B across SQuAD (Exact Match and F1) and Natural Questions (Exact Match). Results are shown for multiple skip strategies and levels. **FiRST (Ours)** consistently performs best under skipping, demonstrating strong robustness. **R + LoRA** indicates Router Augmentation followed by LoRA fine-tuning.

TPOT. At 24%, FiRST is significantly better than other layer skipping strategies, while achieving > 20% improvement in latency.

SQuAD Dataset: For the LLaMA-3-8B model, FiRST (at 10% skip level) maintains over 95% of the performance of the gold standard (Base + LoRA without skipping) (Table 2), with overall latency gains upto 16% (Table 4). It is **significantly** better in quality than all other baselines across all metrics for different levels of skipping. For LLaMA-3.2-3B, again FiRST is > 95% in output quality of gold (base + LoRA) for 10% skipping (Table 2) with similar gains in latency (upto 16% overall) (Table 4) over the LoRA fine-tuned base model. Moreover, it is better than all other layer skipping strategies across all metrics.

Natural Questions Dataset: For the LLaMA-3-

8B model, FiRST retains over 92% of the EM achieved by the non-skipping gold standard (Table 2), with overall latency gains of 5-12% (Table 4). It is **significantly** better in quality than all other baselines for different levels of skipping. For LLaMA-3.2-3B, again FiRST is > 93% in output quality of gold (base + LoRA) for 10% skipping (Table 2) with gains in latency of 4-10% overall (Table 4) over the LoRA fine-tuned base model. Moreover, it is better than all other layer skipping strategies. While more sophisticated retrieval strategies could further enhance performance, our goal was to demonstrate that FiRST maintains quality performance and latency gains even in a retrieval-based QA setting.

Detailed results for an additional skipping percentage are provided in Appendix A.6.

Layer-wise Skipping Patterns: Layer-wise skipping varies significantly across tasks, reflecting the task-specific importance of each layer. For LLaMA-3-8B at a 15% skipping rate, layers 7–9 and 21 are fully skipped in English-to-German, with partial skipping in layer 18. Figure 3 for summarization shows that only some layers in the middle of the network, specifically layers 19 to 27, are skipped. The early layers and the last few layers (28–32) are never skipped. Layers 20, 22, and 23 are fully skipped, with partial skipping in layers 19 and 21. Some layers are skipped less than 10% of times, indicating their necessity for specific sequences. This shows that it’s important to learn which layers to skip based on each input, as the skipping pattern is not the same for every sample. The task-specificity is also evident in Question-Answering, where only layer 22 is fully skipped on SQuAD while layers 12, 23 are fully skipped in case of NQ dataset, and skipping patterns depend on the input. Detailed statistics are in Appendix A.7. Furthermore, the router is trained in a task-aware manner, ensuring that skipping decisions align with task complexity. We also conduct experiments to validate the reusability and generalizability of routers (see Appendix: A.8) across different datasets for the same underlying task.

Computational overhead of Routers: We aim to improve the efficiency of pre-trained models by introducing routers, lightweight linear classifiers that help decide whether to skip certain layers during inference. A common concern with adding such components is the potential increase in computational cost. However, our analysis shows that the

Skip (%)	Model Type		BERT	R-1	R-L
LLaMA-3-8B					
0	Original Model	Base + LoRA	84.87	28.46	16.99
		Base	82.29	23.49	14.66
15	Skip Decode	R + LoRA	84.74	22.04	17.54
		Router	82.53	13.68	9.30
	Random Skip	R + LoRA	83.70	24.60	15.01
		Router	81.10	19.64	13.07
	Unified Skip	R + LoRA	84.25	24.35	14.3
		Router	80.3	16.61	10.95
FiRST (Ours)	R + LoRA	85.14	31.8	20.13	
	Router	81.25	20.2	13.01	
20	Skip Decode	R + LoRA	82.57	20.41	14.87
		Router	81.62	13.48	9.19
	Random Skip	R + LoRA	81.39	21.57	13.83
		Router	79.23	15.51	10.93
	Unified Skip	R + LoRA	82.93	22.3	13.37
		Router	80.32	16.51	11.15
FiRST (Ours)	R + LoRA	82.8	27.65	17.84	
	Router	79.32	16.28	10.85	

Skip (%)	Model Type		BERT	R-1	R-L
LLaMA-3.2-3B					
0	Original Model	Base + LoRA	84.89	28.37	17.02
		Base	71.85	19.34	12.00
15	Skip Decode	R + LoRA	83.20	21.71	13.74
		Router	80.97	9.74	6.87
	Random Skip	R + LoRA	79.52	20.18	12.10
		Router	68.20	10.10	7.10
	Unified Skip	R + LoRA	81.53	18.89	11.72
		Router	70.01	12.49	8.68
FiRST (Ours)	R + LoRA	83.17	26.47	16.79	
	Router	70.98	16.47	10.51	
24	Skip Decode	R + LoRA	78.55	15.83	6.74
		Router	76.91	13.29	8.86
	Random Skip	R + LoRA	80.00	16.33	10.07
		Router	67.88	8.49	5.96
	Unified Skip	R + LoRA	79.31	15.88	10.69
		Router	68.86	9.17	6.97
FiRST (Ours)	R + LoRA	80.25	21.28	13.89	
	Router	69.17	12.36	8.43	

Table 3: Quality Analysis on Summarization (CNN/DM dataset) on LLaMA-3-8B (left) and LLaMA-3.2-3B (right): BERT F1, Rouge-1 and Rouge-L scores are reported for varying skipping levels. Note that R + LoRA corresponds to Router Augmentation followed by LoRA fine-tuning (in the proposed FiRST framework).

Model size	Model Type	WMT			CNN/DM		SQuAD		Natural Questions	
		Skip (%)	Eng→De	Eng→Zh	Skip (%)	TPOT	Skip (%)	TPOT	Skip (%)	TPOT
LLaMA-3-8B	Base + LoRA	0	1×	1×	0	1×	0	1×	0	1×
	R + LoRA	15	0.90×	0.88×	15	0.88×	10	0.95×	10	0.96×
	R + LoRA	25	0.82×	0.83×	20	0.81×	20	0.78×	20	0.80×
LLaMA-3.2-3B	Base + LoRA	0	1×	1×	0	1×	0	1×	0	1×
	R + LoRA	15	0.90×	0.91×	15	0.88×	15	0.94×	15	0.94×
	R + LoRA	25	0.78×	0.75×	24	0.79×	24	0.83×	24	0.84×

Table 4: TPOT variation across all datasets for FiRST. The reported values are relative to the LoRA fine-tuned base model. **R + LoRA** indicates Router Augmentation followed by LoRA fine-tuning. Fine-tuning improves TPOT and quality significantly.

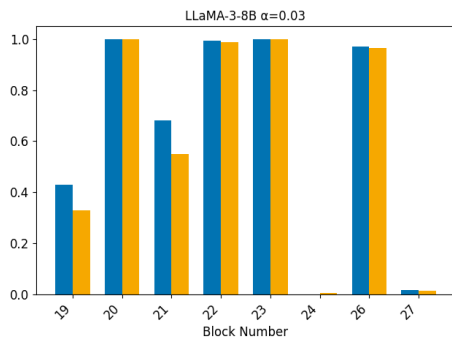


Figure 3: LLaMA-3-8B skipping statistics at 15% skipping rate on summarization task. Layers with no skipping, indicated by a 0% skipping rate, are not represented in the plot.

additional parameters and operations introduced by the routers account for a very small fraction of the total model computation, only 0.0027% for LLaMA-3-8B and 0.0016% for LLaMA-3.2-3B in terms of parameter count, and around 0.3% of the total forward pass time. This fixed cost remains negligible compared to the overall model execution time. As a result, the overall latency gain primarily depends on how much computation is skipped, rather than the routing mechanism itself. Given similar levels of skipping, different methods tend to show comparable efficiency gains, so the key

focus becomes minimizing the performance drop relative to the original model. Our results show that the proposed routing mechanism achieves this balance effectively, preserving latency improvements while maintaining model quality.

6 Conclusion

We propose a new framework, FiRST, for layer selection that adapts to the input sequence and task, aiming to reduce latency in a quality-aware manner. This approach is sequence-dependent and operates compatibly with KV caching. With an optimal layer skipping rate of around 15%, FiRST achieves a 10-20% reduction in latency. This speedup is achieved while remaining quality-neutral, maintaining approximately 80% or more of the base model’s performance on quality metrics across multiple tasks (Machine Translation, Summarization, Question Answering) and model architectures on well-known open-source datasets. Furthermore, our method significantly outperforms other layer selection strategies on most quality metrics.

7 Limitations

FIRST algorithm selects layers in a greedy, myopic way one layer at a time, corresponding to sequences (and tasks). A more optimal way of doing this would be to estimate a subsequence of layers to traverse through instead of one layer at a time. We intend to address this in future work. We would like to select a more optimal subset (subsequence) of layers which will increase the output quality while reducing latency even further.

8 Ethical Concerns

There are no ethical concerns to the best of our knowledge.

References

- Rishabh Agarwal, Nino Vieillard, Piotr Stanczyk, Sabela Ramos, Matthieu Geist, and Olivier Bachem. Gkd: Generalized knowledge distillation for autoregressive sequence models. *arXiv preprint arXiv:2306.13649*, 2023.
- Saleh Ashkboos, Maximilian L Croci, Marcelo Genari do Nascimento, Torsten Hoefler, and James Hensman. Slicept: Compress large language models by deleting rows and columns. *arXiv preprint arXiv:2401.15024*, 2024.
- Xiaodong Chen, Yuxuan Hu, and Jing Zhang. Compressing large language models by streamlining the unimportant layer. *arXiv preprint arXiv:2403.19135*, 2024.
- Luciano Del Corro, Allie Del Giorno, Sahaj Agarwal, Bin Yu, Ahmed Awadallah, and Subhabrata Mukherjee. Skipdecode: Autoregressive skip decoding with batching and caching for efficient llm inference. *arXiv preprint arXiv:2307.02628*, 2023.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024.
- Maha Elbayad, Jiatao Gu, Edouard Grave, and Michael Auli. Depth-adaptive transformer. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=SJg7KhVKPH>.
- Mostafa Elhoushi, Akshat Shrivastava, Diana Liskovich, Basil Hosmer, Bram Wasti, Liangzhen Lai, Anas Mahmoud, Bilge Acun, Saurabh Agarwal, Ahmed Roman, et al. Layer skip: Enabling early exit inference and self-speculative decoding. *arXiv preprint arXiv:2404.16710*, 2024.
- 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*, pp. 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.
- Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. Minillm: Knowledge distillation of large language models. In *The Twelfth International Conference on Learning Representations*, 2024.
- Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. Teaching machines to read and comprehend. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. URL https://proceedings.neurips.cc/paper_files/paper/2015/file/afdec7005cc9f14302cd0474fd0f3c96-Paper.pdf.
- Lu Hou, Zhiqi Huang, Lifeng Shang, Xin Jiang, Xiao Chen, and Qun Liu. Dynabert: Dynamic bert with adaptive width and depth. *Advances in Neural Information Processing Systems*, 33:9782–9793, 2020.
- Ajay Jaiswal, Bodun Hu, Lu Yin, Yeonju Ro, Shiwei Liu, Tianlong Chen, and Aditya Akella. Ffn-skipllm: A hidden gem for autoregressive decoding with adaptive feed forward skipping. *arXiv preprint arXiv:2404.03865*, 2024.
- Wenxiang Jiao, Jen tse Huang, Wenxuan Wang, Zhiwei He, Tian Liang, Xing Wang, Shuming Shi, and Zhaopeng Tu. Parrot: Translating during chat using large language models tuned with human translation and feedback, 2023. URL <https://arxiv.org/abs/2304.02426>.
- Tom Kocmi, R. Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark A. Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, C. Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Mikulas Popovic. Findings of the 2022 conference on machine translation (wmt22). In *Conference on Machine Translation*, pp. 1–45, 2022.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466, 2019. doi: 10.1162/tacl_a_00276. URL <https://aclanthology.org/Q19-1026/>.

- Changhun Lee, Jungyu Jin, Taesu Kim, Hyungjun Kim, and Eunhyeok Park. Owq: Outlier-aware weight quantization for efficient fine-tuning and inference of large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 13355–13364, 2024.
- Lei Li, Yankai Lin, Deli Chen, Shuhuai Ren, Peng Li, Jie Zhou, and Xu Sun. Cascadebert: Accelerating inference of pre-trained language models via calibrated complete models cascade. *arXiv preprint arXiv:2012.14682*, 2020.
- Yixiao Li, Yifan Yu, Chen Liang, Pengcheng He, Nikos Karampatziakis, Weizhu Chen, and Tuo Zhao. Loftq: Lora-fine-tuning-aware quantization for large language models. *arXiv preprint arXiv:2310.08659*, 2023.
- Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Weiming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. Awq: Activation-aware weight quantization for on-device llm compression and acceleration. *Proceedings of Machine Learning and Systems*, 6:87–100, 2024.
- Weijie Liu, Peng Zhou, Zhe Zhao, Zhiruo Wang, Haotang Deng, and Qi Ju. Fastbert: a self-distilling bert with adaptive inference time. *arXiv preprint arXiv:2004.02178*, 2020.
- Yijin Liu, Xianfeng Zeng, Fandong Meng, and Jie Zhou. Instruction position matters in sequence generation with large language models, 2023a. URL <https://arxiv.org/abs/2308.12097>.
- Yijin Liu, Fandong Meng, and Jie Zhou. Accelerating inference in large language models with a unified layer skipping strategy. *arXiv preprint arXiv:2404.06954*, 2024.
- Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang Shi, Raghuraman Krishnamoorthi, and Vikas Chandr. Llm-qat: Data-free quantization aware training for large language models. *arXiv preprint arXiv:2305.17888*, 2023b.
- 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*, pp. 22137–22176. PMLR, 2023c.
- X Ma, G Fang, and X Wang. On the structural pruning of large language models. *NeurIPS, Llm-pruner*, 2023a.
- 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, 2023b.
- Matanel Oren, Michael Hassid, Nir Yarden, Yossi Adi, and Roy Schwartz. Transformers are multi-state rnns. *arXiv preprint arXiv:2401.06104*, 2024.
- Dheeraj Peri, Jhalak Patel, and Josh Park. Deploying quantization-aware trained networks using tensorrt. In *GPU Technology Conference*, 2020.
- Matt Post. A call for clarity in reporting BLEU scores. In Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Christof Monz, Matteo Negri, Aurélie Névéal, Mariana Neves, Matt Post, Lucia Specia, Marco Turchi, and Karin Verspoor (eds.), *Proceedings of the Third Conference on Machine Translation: Research Papers*, pp. 186–191, Brussels, Belgium, October 2018. Association for Computational Linguistics. doi: 10.18653/v1/W18-6319. URL <https://aclanthology.org/W18-6319/>.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100, 000+ questions for machine comprehension of text. *CoRR*, abs/1606.05250, 2016. URL <http://arxiv.org/abs/1606.05250>.
- Rajarshi Saha, Varun Srivastava, and Mert Pilanci. Matrix compression via randomized low rank and low precision factorization. *Advances in Neural Information Processing Systems*, 36, 2023.
- Tal Schuster, Adam Fisch, Jai Gupta, Mostafa Dehghani, Dara Bahri, Vinh Tran, Yi Tay, and Donald Metzler. Confident adaptive language modeling. *Advances in Neural Information Processing Systems*, 35:17456–17472, 2022a.
- Tal Schuster, Adam Fisch, Jai Gupta, Mostafa Dehghani, Dara Bahri, Vinh Q. Tran, Yi Tay, and Donald Metzler. Confident adaptive language modeling, 2022b. URL <https://arxiv.org/abs/2207.07061>.
- Shaohuai Shi, Qiang Wang, and Xiaowen Chu. Efficient sparse-dense matrix-matrix multiplication on gpus using the customized sparse storage format. In *2020 IEEE 26th International Conference on Parallel and Distributed Systems (ICPADS)*, pp. 19–26. IEEE, 2020.
- Jiwon Song, Kyungseok Oh, Taesu Kim, Hyungjun Kim, Yulhwa Kim, and Jae-Joon Kim. Sleb: Streamlining llms through redundancy verification and elimination of transformer blocks. *arXiv preprint arXiv:2402.09025*, 2024.
- Tianxiang Sun, Xiangyang Liu, Wei Zhu, Zhichao Geng, Lingling Wu, Yilong He, Yuan Ni, Guotong Xie, Xuanjing Huang, and Xipeng Qiu. A simple hash-based early exiting approach for language understanding and generation. *arXiv preprint arXiv:2203.01670*, 2022.
- Jiaming Tang, Yilong Zhao, Kan Zhu, Guangxuan Xiao, Baris Kasikci, and Song Han. Quest: Query-aware sparsity for efficient long-context llm inference. *arXiv preprint arXiv:2406.10774*, 2024.
- Ziheng Wang. Sparsert: Accelerating unstructured sparsity on gpus for deep learning inference. In *Proceedings of the ACM international conference on parallel*

architectures and compilation techniques, pp. 31–42, 2020.

- Chris Wendler, Veniamin Veselovsky, Giovanni Monea, and Robert West. Do llamas work in English? on the latent language of multilingual transformers. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15366–15394, 2024. URL <https://aclanthology.org/2024.acl-long.820>.
- Yuxiang Wu, Sebastian Riedel, Pasquale Minervini, and Pontus Stenetorp. Don’t read too much into it: Adaptive computation for open-domain question answering, 2020. URL <https://arxiv.org/abs/2011.05435>.
- Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming language models with attention sinks. *arXiv preprint arXiv:2309.17453*, 2023.
- Ji Xin, Raphael Tang, Jaejun Lee, Yaoliang Yu, and Jimmy Lin. Deebert: Dynamic early exiting for accelerating bert inference. *arXiv preprint arXiv:2004.12993*, 2020.
- Linfeng Zhang, Zhanhong Tan, Jiebo Song, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Scan: A scalable neural networks framework towards compact and efficient models. *Advances in Neural Information Processing Systems*, 32, 2019.
- Yuxin Zhang, Lirui Zhao, Mingbao Lin, Yunyun Sun, Yiwu Yao, Xingjia Han, Jared Tanner, Shiwei Liu, and Rongrong Ji. Dynamic sparse no training: Training-free fine-tuning for sparse llms. *arXiv preprint arXiv:2310.08915*, 2023.
- Wangchunshu Zhou, Canwen Xu, Tao Ge, Julian McAuley, Ke Xu, and Furu Wei. Bert loses patience: Fast and robust inference with early exit. *Advances in Neural Information Processing Systems*, 33:18330–18341, 2020.
- Wei Zhu. Leebert: Learned early exit for bert with cross-level optimization. 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)*, pp. 2968–2980, 2021.
- Xunyu Zhu, Jian Li, Yong Liu, Can Ma, and Weiping Wang. A survey on model compression for large language models. *arXiv preprint arXiv:2308.07633*, 2023.

A Appendix

A.1 Related Work

Early Exit: Several works have been proposed in the early exit theme (Zhu, 2021; Zhou et al., 2020; Xin et al., 2020; Liu et al., 2020; Li et al., 2020; Hou et al., 2020; Schuster et al., 2022a; Wu et al., 2020) where adaptive compute is used for different parts of the token sequence. While these approaches have been popular for encoder-only models which processes the entire sequence as a whole, they have faced challenges in generation tasks. The main limitation of these set of techniques are their inability to handle KV caching appropriately which is crucial for multi-fold speed-ups in current LLM architectures. We emphasize that in our work, we assign varying compute to sequences in different batches but within the same sequence, we assign the same compute to every token.

Input Agnostic Heuristics: In Skip Decoding (Del Corro et al., 2023), initial tokens pass through more layers than later ones, contradicting the observation that later tokens are harder to decode (Liu et al., 2024). Additionally, Skip Decoding skips several bottom layers for most tokens, causing undesirable sub-network imbalance. To address this, Unified Layer Skipping (Liu et al., 2024) proposes a discrete skipping strategy that is uniform for all tokens in a sequence. Based on a latency budget, retained layer IDs are passed through by all tokens, ensuring KV Cache handling and retaining key layers. However, the limitation of this approach is that skipping is independent of the input sequence. In contrast, early exit strategies adapt layer skipping to the input sequence, offering more flexibility. In (Fan et al., 2019), a method akin to dropout randomly skips layers during training, but this leads to performance decline during the pre-fill stage. FFN-SkipLLM (Jaiswal et al., 2024) constrains skipping to FFN layers to avoid KV Cache issues, but fails to fully exploit redundancy as discussed already. (Song et al., 2024) is a very recent work that also explores greedily identifying layers to skip while preserving the model performance on a calibration dataset - however there are two major limitations of this work which are resolved in our paper: (A) first of all, the layer selection strategy is sequence independent although it can be made task-dependent by calibrating on task-specific datapoints, our approach for skipping layers is sequence dependent and is based on the input to a layer (B) SLEB does

not explore the impact of fine-tuning on the layers to be skipped. On the other hand, our skipping strategy incorporates the trained router already - intuitively, the knowledge of skipping is transferred during finetuning. (Chen et al., 2024) is another recent work that compresses models by identifying redundant layers - this is done by computing the average similarity between input/output pairs of a layer. However, as outlined in (Song et al., 2024), such an approach suffers from the limitation that it does not take into account the joint association between the layers while skipping multiple layers. Moreover, like (Jaiswal et al., 2024), this work is neither sequence dependent nor takes the final model predictions into account while identifying the layers to skip.

Model Compression and Quantization Aware Training:

Orthogonal approaches to explore the latency/memory-performance trade-off in Large Language Models aim to build smaller models that approximate the performance of larger ones with reduced memory and latency costs. Key techniques include: 1) compressing model parameters into fewer bits (Frantar et al., 2022; Lin et al., 2024; Lee et al., 2024; Saha et al., 2023); 2) pruning the network by removing components like attention heads or neurons based on heuristics (Frantar & Alistarh, 2023; Ma et al., 2023b); and 3) distilling the large model into a smaller, faster counterpart (Agarwal et al., 2023; Gu et al., 2024). For further details, we refer to the survey by (Zhu et al., 2023). A significant body of work (Dettmers et al., 2024; Liu et al., 2023b; Peri et al., 2020; Li et al., 2023) has focused on quantization-aware training to reduce memory footprints and mitigate performance loss, starting with QLoRA (Dettmers et al., 2024). In a similar vein, our work proposes finetuning router-augmented models to improve layer skipping and reduce performance degradation, as pre-trained models do not account for layer skipping, leading to higher degradation with vanilla skipping.

Network Pruning: Another orthogonal approach to improve the inference speed-up is to prune redundant network weights by zeroing them out. There has been a significant body of work on pruning model weights (Frantar et al., 2022; Frantar & Alistarh, 2023; Sun et al., 2022; Zhang et al., 2023) - most of these works can be categorized into two clusters namely unstructured pruning and

structured pruning. In case of unstructured pruning, there is no structure to the inserted zeros and achieving speedups with modern GPU hardware tailored towards dense matrix multiplication is challenging. In fact, more than 90% sparsity is typically required to achieve any significant speedup (Wang, 2020; Shi et al., 2020). Therefore, structured pruning which is more amenable to GPU hardware has become prominent (2:4 pruning and sub-channel pruning). However, realizing desired speedups through these techniques have been difficult (Song et al., 2024). Moreover, several approaches for dynamically deleting entire rows or columns of weight matrices have been proposed (Ma et al., 2023a; Ashkboos et al., 2024; Liu et al., 2023c) to retain dense matrices but two limitations remain - (A) hardware support is extremely limited for realizing speedup gains (B) extensive finetuning is necessary to align the sparsification with linguistic abilities - this is because, such pruning techniques were not observed by the model during pre-training. Finally, note that several prior works (Tang et al., 2024; Oren et al., 2024; Xiao et al., 2023) have imposed (query aware/ query agnostic) sparsity in the KV cache matrices to speed up self-attention mechanism via clever selection of the critical tokens necessary from the KV cache.

Mixture of Experts Mixture of Experts (MoE) is a well-established technique for improving the efficiency and capacity of deep learning models by conditionally activating subsets of parameters for different inputs. MoE-based transformer models, such as Switch Transformers (Fedus et al., 2021) and GShard (Lepikhin et al., 2020), employ a gating mechanism to route tokens to a subset of expert layers, thereby significantly reducing computational costs while maintaining expressivity. These architectures are designed to scale up model capacity without a proportional increase in inference cost. The main goal is to scale up the parameters while maintaining the cost of pre-training and inference.

Although MoE and input-adaptive layer skipping share the goal of selective computation, they differ in fundamental ways. MoE dynamically routes tokens to different experts at the layer level, whereas layer skipping focuses on bypassing entire layers in the transformer stack based on the input query. MoE models typically maintain a full-depth model structure, leveraging sparse activation to reduce computational overhead, whereas layer-skipping methods explicitly modify the depth of computa-

tion for different inputs. Simply put, the techniques are orthogonal - layer skipping can be used in tandem with Moe to further reduce the depth of computation wherever possible at an expert level

A.2 Details of Datasets

Split	WMT		Summarization	Question Answering	
	En→De	En→Zh	CNN/DM	SQuAD	NQ
Train	3,505	8,983	3,400	3,400	3,400
Validation	876	998	600	600	600
Test	2,038	2,038	11,490	10,570	7,830

Table 5: Train, validation, and test splits for machine translation (WMT), summarization (CNN/DM), and question answering (SQuAD, NQ) tasks.

Machine Translation: For translation tasks, namely English-to-Chinese and English-to-German, we employ the WMT development sets from 2017 to 2020 for training/fine-tuning following the methodology outlined in previous studies (Liu et al., 2023a; Jiao et al., 2023). Translation performance is evaluated using the test set from the WMT 2022 dataset (Kocmi et al., 2022) which was developed using recent content from diverse domains. These domains include news, social media, e-commerce, and conversational contexts.

Summarization: We use the popular CNN-DailyMail (CNN/DM) (Hermann et al., 2015) dataset which is a large collection (over 300k) of text summarization pairs, created from CNN and Daily Mail news articles. Each datapoint in this dataset comprises of an article (the body of the news article with 683 words on average) and the corresponding highlights (article summary as written by the article author). While the training set contains more than 287k samples, we have randomly chosen 4k samples for training both routers and LoRA. During training in our framework, the number of trainable parameters is small in both phases - therefore a small subset of data points is sufficient for training.

Question Answering: We use the popular **Stanford Question Answering Dataset (SQuAD v1.1)** (Rajpurkar et al., 2016), a widely-used benchmark for Machine Question Answering. The dataset consists of over 100k question-answer pairs posed by crowd-workers on a set of over 500 Wikipedia articles. Each sample comprises a context (a passage from a Wikipedia article), a question (crafted to test comprehension of the passage), and the corresponding answer (a text span from the correspond-

ing reading passage). Similarly to the CNN/DM dataset, 4k samples are chosen at random to train both routers and LoRA. The training and validation splits contain 87,599 and 10,570 samples, respectively. Evaluation is performed on the validation set (Schuster et al., 2022a) as the test set labels are not publicly released. **Natural Questions (NQ):** We use the Natural Questions dataset (Kwiatkowski et al., 2019), which consists of real queries issued to Google Search paired with relevant Wikipedia articles. Each example in NQ contains a query (an actual user question), a context (Wikipedia article), and two types of answers: a long answer (typically a paragraph) and a short answer (a specific text span, when available). We implement a naïve RAG solution using a simple S-BERT model to generate embeddings for the query and passages and retrieve the top 5 most relevant passages based on similarity. Once retrieved, these passages are used as context to answer the query and then compared against gold answers. Similar to SQuAD, 4k samples are randomly chosen to train both routers and LoRA. The training set contains over 300k examples, with evaluation performed on 7.83k validation samples containing short answers. For further details on training-testing split, refer to table: 5

A.3 Evaluation Metrics

Quality-Based Metrics for Translation task:

- **BLEU Score:** BLEU (Bilingual Evaluation Understudy) scores are used to measure the quality of translations. BLEU compares n-grams of the candidate translation to n-grams of the reference translation, providing a score between 0 and 1, with higher scores indicating better translations. In this evaluation, NLTK BLEU is employed and We report BLEU-1, BLEU-2, and the cumulative BLEU score, which is computed as the geometric mean of individual n-gram precision scores from unigram to 4-gram. To mitigate the issue of zero counts for higher-order n-grams, we apply the smoothing strategy utilized in (Post, 2018).
- **COMET:** COMET (Cross-lingual Optimized Metric for Evaluation of Translation) is used to assess translation quality further. COMET evaluates translations using a model trained to correlate well with human judgments. Specifically, Unbabel/XCOMET-XL² is used in this evaluation. COMET provides a more nuanced assessment of translation quality by considering the

²<https://github.com/Unbabel/COMET>

intricacies of both source and target languages, beyond the n-gram matching used in BLEU.

Quality based Metrics for Summarization Task:

- **BERTScore:** This metric quantifies semantic similarity between texts by leveraging contextual word embeddings. BERTScore captures meaning-based similarity rather than relying on exact word matches, providing a nuanced evaluation of text generation quality.
- **ROUGE:** (Recall-Oriented Understudy for Gisting Evaluation) is a common metric - ROUGE-1 refers to overlap of unigrams between the system summary and reference summary. Similarly, ROUGE-L measures longest matching sequence of words.

Quality based Metrics for Question Answering Task:

- **Exact Match:** This metric measures the percentage of predictions that exactly match the ground truth answer.
- **F1 score:** Since EM is a highly stringent metric, we also report the F1 score which provides a more flexible evaluation of answer prediction. This metric also takes into account near-matches.

A.4 Training and Inference Setup

- **Training settings:** We perform extensive experiments on two models, namely LLaMA-3-8B and LLaMA-3.2-3B from Meta, which consist of 32 and 28 layers, respectively. Training of routers and LoRA adapters is conducted on A100 80GB GPUs, with training/inference is performed in full precision to avoid performance degradation due to quantization. The training process employs our custom loss function and continues for a fixed number of epochs, terminating when the validation loss fails to improve over 4 consecutive steps. The learning rate is set between $1e^{-4}$ and $3e^{-4}$ - a cosine scheduler is used to adjust the learning rate. Gradients are accumulated after five steps, and we set the regularization loss coefficient ($\lambda = 0.01$), ensuring it meaningfully contributes to the overall loss without overpowering primary losses like cross-entropy and penalization loss. After training, we verify the router weight norms (e.g., thresholds of 0.1 or 0.05) to ensure they remain stable, neither exploding nor vanishing, preventing overfitting or underfitting. For LoRA fine-tuning, we employ a rank

of 8, a dropout rate of 0.1, and a scaling factor ($\text{lor}\alpha = 32$). Our approach involves two key penalization coefficients: α for the router and β for LoRA training. The non-skip coefficient (α) is adjusted based on dataset characteristics, sequence length, and model depth, requiring some tuning. Typically, router penalization (α) is set $2-4\times$ higher than LoRA's (β) based on experimental observations to maintain a consistent skipping percentage across both training phases. A reasonable starting point is to tune α for a 10-15% skipping range, then gradually increase it at regular intervals to encourage higher skipping levels. In addition to the above hyperparameters, we define the maximum sequence length based on the task. For translation, it is set to 128 for router training and 256 for LoRA training. Similarly, for summarization, the sequence length is set to 500 and 700, respectively. For Question Answering, this length is set to 512. The prompts for the different tasks related to training / inference are shown in Appendix A.5.

- **Joint training considerations:** The decision to split training into two phases, first training the router and then the LoRA modules, was made to ensure better stability during optimization. In our experiments, this approach allowed the router to establish a robust skipping strategy without interference from the fine-tuning of LoRA modules, which otherwise introduces instability. While training both the router and LoRA modules together is theoretically possible and could eliminate the need for reapplying the non-skip penalty in the second phase, initial experiments revealed that it led to suboptimal convergence and conflicting gradients. This was particularly evident when the router's decisions were not yet stable, causing inconsistencies in the joint optimization process.
- **Inference settings:** For all the tasks, we set the temperature to 0.8 and enable top-k sampling over 10 tokens. The maximum number of tokens to be generated is set to 80 for WMT, 200 for CNN/DM and 32 for SQuAD and NQ. Caching is turned on during inference.

A.5 Prompt Details

The prompt structures used for both training and inference are as follows:

- For the machine translation task (English-to-

German or English-to-Chinese), the following general prompt structure is used to train the routers and during final inference:

```

### Instruction:
Translate the following sentences from English to German.

### Input:
{Text to be translated}

### Response:

```

- For the summarization task (used in CNN/DailyMail dataset), the prompt structure used is:

```

### Instruction:
Summarize the news article in around 100-200 words.

### Input:
{Article to be summarized}

### Response:

```

- For the Question Answering task (used in SQuAD/NQ dataset), the following prompt structure is utilized:

```

### Instruction:
Answer the question based on the given passage.

### Passage:
{context}

### Question:
{Question to be answered}

### Response:

```

During the training of the LoRA module, task-aware training is applied. The expected translation or summary is appended after the `### Response` section, making the model predict the response tokens following the "Response:\n".

A.6 Detailed Result Table

Tables 6 and 7 present the detailed results for the LLaMA-3.2-3B and LLaMA-3-8B models for an additional skipping percentage on Machine Translation Task. The results are reported using BLEU (BLEU-1, BLEU-2, BLEU) and COMET metrics, highlighting performance across different skipping percentages. Similarly, Table 11 presents cumulative results for both models reporting BERT F1, ROUGE-1 and ROUGE-L. Lastly, Table 8 presents Exact Match (EM) and F1 scores for both models for three skipping percentage variations.

Skip (%)	Model Type	Config	BLEU-1	BLEU-2	BLEU	COMET
LLaMA-3.2-3B						
0	Original Model	Base + LoRA	0.376	0.174	0.160	89.72
		Base	0.312	0.137	0.125	81.66
15	Skip Decode	Router + LoRA	0.207	0.069	0.105	44.58
		Router	0.099	0.039	0.055	32.74
	Random Skip	Router + LoRA	0.190	0.048	0.079	47.26
		Router	0.121	0.028	0.048	36.30
	Unified Skip	Router + LoRA	0.216	0.059	0.095	44.72
		Router	0.164	0.040	0.067	39.81
	FiRST (Ours)	Router + LoRA	0.247	0.071	0.113	60.29
		Router	0.165	0.041	0.069	43.04
25	Skip Decode	Router + LoRA	0.168	0.047	0.077	32.33
		Router	0.091	0.035	0.050	27.64
	Random Skip	Router + LoRA	0.090	0.011	0.023	30.97
		Router	0.053	0.007	0.015	27.22
	Unified Skip	Router + LoRA	0.151	0.027	0.050	39.81
		Router	0.102	0.016	0.032	30.86
	FiRST (Ours)	Router + LoRA	0.189	0.041	0.072	45.38
		Router	0.095	0.014	0.029	29.08
35	Skip Decode	Router + LoRA	0.058	0.011	0.018	23.30
		Router	0.017	0.003	0.005	19.34
	Random Skip	Router + LoRA	0.047	0.001	0.007	25.08
		Router	0.020	0.001	0.003	21.18
	Unified Skip	Router + LoRA	0.012	0.000	0.002	19.65
		Router	0.009	0.000	0.001	20.36
	FiRST (Ours)	Router + LoRA	0.095	0.013	0.027	27.45
		Router	0.057	0.006	0.014	25.03
LLaMA-3-8B						
0	Original Model	Base + LoRA	0.418	0.217	0.199	93.00
		Base	0.372	0.186	0.169	87.13
15	Skip Decode	Router + LoRA	0.230	0.105	0.094	55.62
		Router	0.040	0.012	0.019	23.33
	Random Skip	Router + LoRA	0.304	0.110	0.097	66.25
		Router	0.265	0.088	0.132	60.27
	Unified Skip	Router + LoRA	0.289	0.106	0.153	59.34
		Router	0.232	0.079	0.117	59.26
	FiRST (Ours)	Router + LoRA	0.380	0.179	0.161	82.14
		Router	0.288	0.118	0.108	67.74
25	Skip Decode	Router + LoRA	0.118	0.052	0.070	31.47
		Router	0.032	0.009	0.015	21.55
	Random Skip	Router + LoRA	0.060	0.009	0.018	29.71
		Router	0.037	0.005	0.011	29.95
	Unified Skip	Router + LoRA	0.157	0.034	0.060	31.69
		Router	0.126	0.027	0.048	32.15
	FiRST (Ours)	Router + LoRA	0.178	0.041	0.071	34.95
		Router	0.097	0.014	0.029	26.01
35	Skip Decode	Router + LoRA	0.049	0.020	0.028	23.85
		Router	0.030	0.008	0.013	20.03
	Random Skip	Router + LoRA	0.018	0.001	0.004	25.56
		Router	0.014	0.001	0.002	25.34
	Unified Skip	Router + LoRA	0.064	0.008	0.017	22.05
		Router	0.039	0.005	0.011	22.88
	FiRST (Ours)	Router + LoRA	0.064	0.004	0.012	19.96
		Router	0.037	0.001	0.005	21.41

Table 6: Machine Translation Results for English to German on LLaMA-3.2-3B and LLaMA-3-8B: BLEU and COMET scores for various skipping strategies.

A.7 Layer-wise Skipping Statistics

To illustrate broader layer-wise behaviors, Figures 4 through 8 show block-wise skip rates observed when the models were configured for an average skip rate of approximately 15%. These figures highlight how individual layers exhibit distinct skipping behaviors across various tasks, including translation (English-German, English-Chinese), summarization (CNN/DM), and question answering (SQuAD, Natural Questions). It also reveals that tasks with similar processing requirements (e.g., different translation language pairs) exhibit comparable skipping patterns. Providing a more granular view for specific models, Table 9 reports the fraction of sequences that skip each block un-

Skip (%)	Model Type	Config	BLEU-1	BLEU-2	BLEU	COMET
LLaMA-3.2-3B						
0	Original Model	Base + LoRA	0.518	0.300	0.278	79.13
		Base	0.321	0.179	0.166	61.84
15	Skip Decode	Router + LoRA	0.381	0.217	0.200	46.70
		Router	0.096	0.048	0.062	34.14
	Random Skip	Router + LoRA	0.387	0.174	0.154	57.79
		Router	0.136	0.056	0.077	35.86
	Unified Skip	Router + LoRA	0.370	0.173	0.157	57.10
		Router	0.224	0.094	0.087	45.16
	FiRST (Ours)	Router + LoRA	0.457	0.237	0.218	67.45
		Router	0.227	0.109	0.100	54.55
25	Skip Decode	Router + LoRA	0.278	0.157	0.147	42.01
		Router	0.070	0.034	0.045	29.51
	Random Skip	Router + LoRA	0.251	0.099	0.137	45.53
		Router	0.101	0.037	0.053	31.50
	Unified Skip	Router + LoRA	0.306	0.123	0.108	42.30
		Router	0.152	0.055	0.079	31.63
	FiRST (Ours)	Router + LoRA	0.329	0.137	0.126	41.66
		Router	0.105	0.035	0.053	27.83
35	Skip Decode	Router + LoRA	0.028	0.011	0.016	22.30
		Router	0.024	0.010	0.014	18.32
	Random Skip	Router + LoRA	0.042	0.015	0.022	32.81
		Router	0.021	0.008	0.011	23.63
	Unified Skip	Router + LoRA	0.033	0.005	0.010	19.12
		Router	0.022	0.004	0.007	19.82
	FiRST (Ours)	Router + LoRA	0.100	0.070	0.075	28.10
		Router	0.065	0.018	0.029	23.13
LLaMA-3-8B						
0	Original Model	Base + LoRA	0.569	0.356	0.333	82.66
		Base	0.380	0.225	0.208	68.95
15	Skip Decode	Router + LoRA	0.287	0.158	0.149	55.98
		Router	0.047	0.023	0.030	21.75
	Random Skip	Router + LoRA	0.479	0.258	0.237	67.32
		Router	0.366	0.187	0.168	59.89
	Unified Skip	Router + LoRA	0.466	0.250	0.230	69.58
		Router	0.273	0.134	0.122	54.57
	FiRST (Ours)	Router + LoRA	0.484	0.266	0.247	68.63
		Router	0.176	0.087	0.080	42.76
25	Skip Decode	Router + LoRA	0.173	0.094	0.088	33.85
		Router	0.038	0.018	0.024	20.93
	Random Skip	Router + LoRA	0.117	0.047	0.065	27.73
		Router	0.074	0.028	0.040	35.16
	Unified Skip	Router + LoRA	0.349	0.158	0.142	50.59
		Router	0.177	0.074	0.068	38.74
	FiRST (Ours)	Router + LoRA	0.358	0.157	0.119	56.92
		Router	0.110	0.032	0.051	25.45
35	Skip Decode	Router + LoRA	0.112	0.057	0.054	25.23
		Router	0.036	0.017	0.023	22.84
	Random Skip	Router + LoRA	0.074	0.022	0.034	27.35
		Router	0.045	0.011	0.018	29.46
	Unified Skip	Router + LoRA	0.075	0.021	0.034	20.25
		Router	0.039	0.011	0.017	21.24
	FiRST (Ours)	Router + LoRA	0.157	0.040	0.066	26.80
		Router	0.061	0.015	0.025	22.89

Table 7: Machine Translation Results for English to Chinese on LLaMA-3.2-3B and LLaMA-3-8B: BLEU and COMET scores for various skipping strategies.

der LoRA adaptation for the LLaMA-3-8B and LLaMA-3.2-3B models. These fractions are presented as decimals, so a value of 0.10 signifies that 10% of sequences skip the corresponding block.

A.8 Generalizability and Reusability of Routers

We also investigate whether a router trained on one dataset can be applied directly to another dataset for the same task, without any additional training. Analysis of layer-skipping patterns suggests that models processing different datasets of the same task exhibit similar internal behavior, implying that a router learned on one dataset could be reused elsewhere. To validate this idea, we train routers on

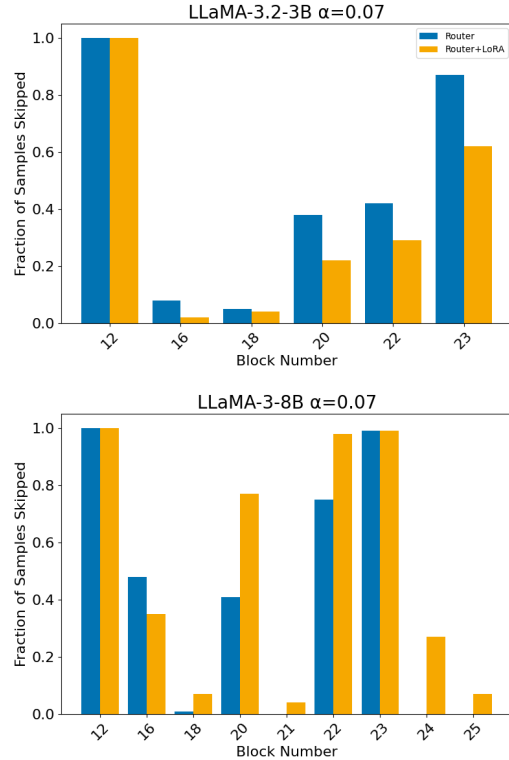


Figure 4: Comparison of LLaMA-3.2-3B (top) and LLaMA-3-8B (bottom) at 10% skipping rate on the Natural Questions Answering Task. Layers with no skipping, indicated by a 0% skipping rate, are not represented in the plot.

English→German translation data and test them on English→Chinese translation, and likewise train on SQuAD for question answering before evaluating on the Natural Questions dataset. The quantitative results of these cross-dataset evaluations are presented in Table 10.

In the QA experiments, the router trained and tested on NQ itself (denoted as *Original* in Table 10) and the router trained on SQuAD but tested on NQ (*Cross-dataset*) both skip a comparable fraction of layers. Despite being trained on different data, the SQuAD-trained router achieves QA accuracy (EM score) very close to that of the NQ-trained router, as detailed in Table 10. This indicates that the essential decision patterns learned by the router transfer well across QA datasets when using a similar level of layer skipping. A similar pattern emerges in machine translation. Whether trained on English→German or directly on English→Chinese data, routers that skip the same proportion of layers produce very similar translation quality metrics (BLEU and COMET in Table 10) on the English→Chinese task. This holds across different model sizes and skip settings. Together, these results demonstrate that routers learned on

one dataset can be effectively reused on another dataset for the same task, as long as they are configured to skip a similar amount of computation. This reusability can lead to substantial savings in both training time and computing resources.

Skip (%)	Model Type	Config	EM	F1
LLaMA-3-8B				
0	Original Model	wLoRA Base	73.93	85.99
			19.46	36.73
10	Skip Decode	R + LoRA Router	60.14	65.33
			16.38	31.48
	Random Skip	R + LoRA Router	65.73	80.08
			18.25	33.75
	Unified Skip	R + LoRA Router	55.54	74.58
			17.39	32.91
FiRST (Ours)	R + LoRA Router	70.85	83.61	
		14.58	31.52	
20	Skip Decode	R + LoRA Router	45.00	55.10
			10.68	26.69
	Random Skip	R + LoRA Router	47.79	66.37
			6.71	22.46
	Unified Skip	R + LoRA Router	52.87	69.28
			18.18	32.51
FiRST (Ours)	R + LoRA Router	60.60	75.49	
		13.21	27.48	
30	Skip Decode	R + LoRA Router	30.77	48.38
			10.67	28.52
	Random Skip	R + LoRA Router	25.45	42.68
			3.55	15.49
	Unified Skip	R + LoRA Router	25.61	38.55
			15.19	28.11
FiRST (Ours)	R + LoRA Router	38.20	52.68	
		3.64	13.30	
LLaMA-3.2-3B				
0	Original Model	wLoRA Base	73.07	84.17
			18.92	37.74
10	Skip Decode	R + LoRA Router	60.79	75.00
			20.00	31.55
	Random Skip	R + LoRA Router	64.78	77.27
			13.76	28.59
	Unified Skip	R + LoRA Router	65.03	77.53
			13.16	32.31
FiRST (Ours)	R + LoRA Router	69.44	81.35	
		12.79	28.37	
20	Skip Decode	R + LoRA Router	40.12	40.00
			20.45	37.62
	Random Skip	R + LoRA Router	11.32	38.34
			6.75	15.51
	Unified Skip	R + LoRA Router	37.39	52.49
			7.81	18.20
FiRST (Ours)	R + LoRA Router	39.70	54.59	
		5.52	15.33	
30	Skip Decode	R + LoRA Router	20.68	23.08
			0.78	3.83
	Random Skip	R + LoRA Router	0.30	8.87
			0.62	5.88
	Unified Skip	R + LoRA Router	7.39	13.72
			0.37	6.15
FiRST (Ours)	R + LoRA Router	33.99	50.37	
		2.55	10.26	

Table 8: SQuAD performance on LLaMA-3-8B and LLaMA-3.2-3B: EM (Exact Match) and F1 scores for varying skipping levels. R + LoRA = Router augmentation + LoRA fine-tuning (FiRST framework). wLoRA = Base model with LoRA fine-tuning.

Layer ↓	$\alpha = 0.005$		$\alpha = 0.01$		$\alpha = 0.025$	
	R	R+L	R	R+L	R	R+L
0	0.0	0.0	0.0	0.0	0.0	0.0
1	0.0	0.0	0.0	0.0	0.0	0.0
2	0.0	0.0	0.0	0.0	0.0	0.0
3	0.0	0.0	0.0	0.0	0.0	0.0
4	0.0	0.0	0.0	0.0	0.0	0.0
5	0.0	0.0	0.0	0.0	0.0	0.0
6	0.0	0.0	0.0	0.0	0.0	0.0
7	100.0	100.0	100.0	100.0	100.0	100.0
8	100.0	100.0	100.0	100.0	100.0	100.0
9	100.0	100.0	0.0	0.0	0.1	1.3
10	0.0	0.0	0.0	0.0	89.3	72.7
11	0.0	0.0	0.0	0.0	0.0	0.0
12	0.0	0.0	100.0	100.0	100.0	100.0
13	0.0	0.0	0.0	0.0	0.0	0.0
14	0.0	0.0	0.0	0.6	0.0	0.3
15	0.0	0.0	100.0	100.0	100.0	99.4
16	0.0	0.0	0.1	2.5	100.0	100.0
17	0.0	0.0	0.0	0.0	0.0	0.0
18	97.8	32.2	100.0	100.0	100.0	100.0
19	0.0	0.0	99.9	91.2	100.0	99.6
20	0.0	0.0	100.0	100.0	100.0	99.5
21	99.9	98.7	100.0	100.0	100.0	100.0
22	0.0	0.0	0.0	0.0	0.0	0.0
23	0.0	0.0	24.1	0.8	99.8	91.7
24	0.0	0.0	0.0	0.0	0.0	0.0
25	0.0	0.0	0.0	0.0	0.0	0.0
26	0.0	0.0	57.3	2.5	100.0	86.0
27	0.0	0.0	0.0	0.0	0.0	0.0
28	0.0	0.0	0.0	0.0	0.0	0.1
29	0.0	0.0	0.0	0.0	0.0	0.0
30	0.0	0.0	0.0	0.0	0.0	0.0
31	0.0	0.0	0.0	0.0	0.0	0.0
Avg	15.6	13.5	27.5	24.9	37.2	36.0

Table 9: Variation in skipping percentage (15–35%) with the non-skip penalization loss coefficient α for LLaMA-3-8B on English–German translation. As α increases, average skipping rises across both Router-only (R) and Router+LoRA (R+L) models.

Question Answering (Natural Questions)			
Setting	Skip (%)	EM	
LLaMA-3.2-3B			
Original	10.32	22.55	
Cross-Dataset	12.46	20.87	
Original	18.03	17.51	
Cross-Dataset	19.04	16.93	
LLaMA-3-8B			
Original	11.44	33.53	
Cross-Dataset	13.11	29.70	
Original	24.60	18.02	
Cross-Dataset	26.80	19.58	
Machine Translation (Eng→Zh)			
Setting	Skip (%)	BLEU	COMET
LLaMA-3.2-3B			
Original	14.69	0.37	0.55
Cross-Dataset	16.57	0.37	0.45
Original	26.12	0.25	0.28
Cross-Dataset	26.16	0.25	0.30
LLaMA-3-8B			
Original	16.46	0.25	0.43
Cross-Dataset	14.70	0.37	0.49
Original	27.73	0.25	0.25
Cross-Dataset	30.00	0.25	0.24

Table 10: **Top**: Exact Match (EM) on Natural Questions (NQ) with two skip settings: *Original*: routers trained/tested on NQ; *Cross-dataset*: trained on SQuAD, tested on NQ. **Bottom**: BLEU and COMET for En→Zh translation: *Original*: trained/tested on En→Zh; *Cross-dataset*: trained on En→De, tested on En→Zh. These results assess router reusability and generalizability across datasets within the same task.

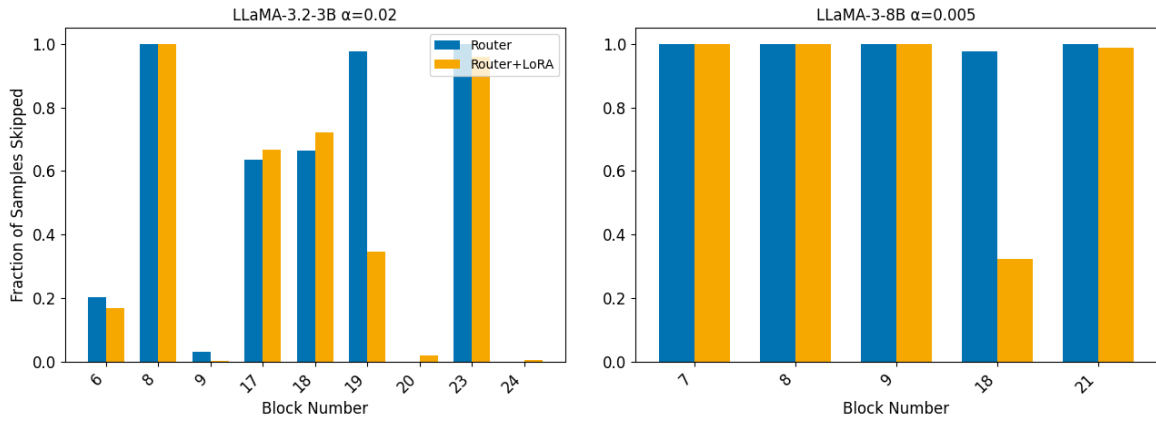


Figure 5: Comparison of LLaMA-3.2-3B (left) and LLaMA-3-8B (right) at 15% skipping rate on English-to-German Machine Translation Task. The graph shows how different layers contribute to the skipping behavior for the same dataset. Layers with no skipping, indicated by a 0% skipping rate, are not represented in the plot.

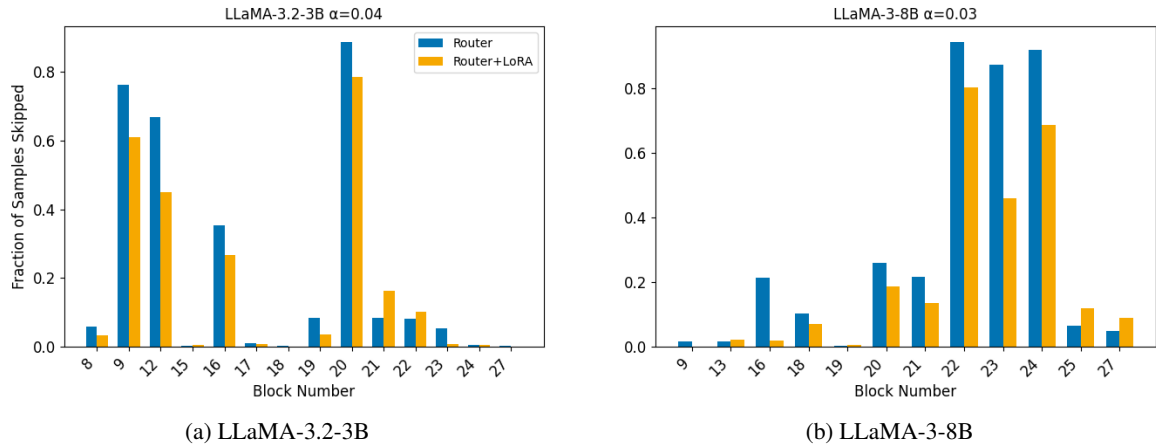


Figure 6: Comparison of LLaMA-3.2-3B (left) and LLaMA-3-8B (right) at 15% skipping rate on SQuAD Question-Answering Task. Layers with no skipping, indicated by a 0% skipping rate, are not represented in the plot.

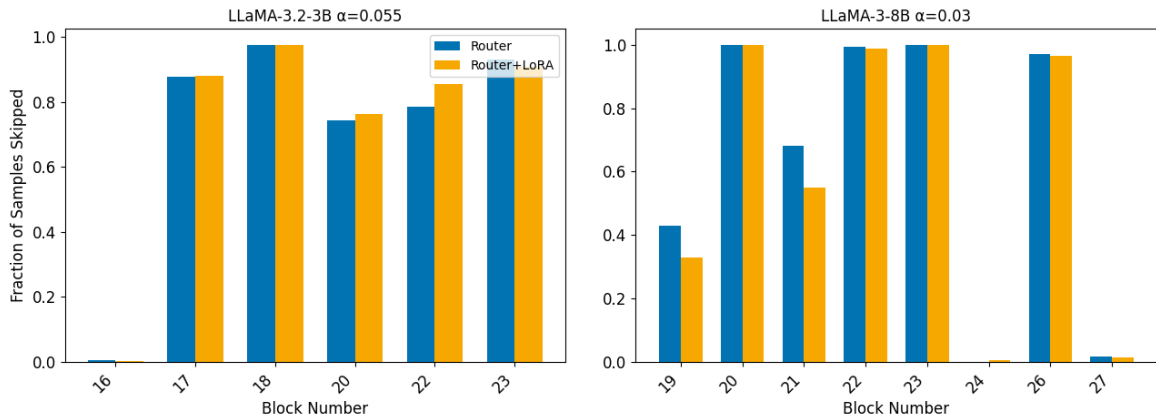


Figure 7: Comparison of LLaMA-3.2-3B (left) and LLaMA-3-8B (right) at 15% skipping rate on CNN Summarization Task. Layers with no skipping, indicated by a 0% skipping rate, are not represented in the plot.

Skip (%)	Model Type		BERT	R-1	R-L
0	Original Model	wLoRA	84.87	28.46	16.99
		Base	82.29	23.49	14.66
15	Skip Decode	R + LoRA	84.74	22.04	17.54
		Router	82.53	13.68	9.30
	Random Skip	R + LoRA	83.70	24.60	15.01
		Router	81.10	19.64	13.07
	Unified Skip	R + LoRA	84.25	24.35	14.3
		Router	80.3	16.61	10.95
FiRST (Ours)	R + LoRA	85.14	31.8	20.13	
Router		81.25	20.2	13.01	
20	Skip Decode	R + LoRA	82.57	20.41	14.87
		Router	81.62	13.48	9.19
	Random Skip	R + LoRA	81.39	21.57	13.83
		Router	79.23	15.51	10.93
	Unified Skip	R + LoRA	82.93	22.3	13.37
		Router	80.32	16.51	11.15
FiRST (Ours)	R + LoRA	82.8	27.65	17.84	
Router		79.32	16.28	10.85	
27	Skip Decode	R + LoRA	79.92	10.67	10.32
		Router	77.27	9.59	7.00
	Random Skip	R + LoRA	76.40	11.45	7.89
		Router	77.45	12.56	9.08
	Unified Skip	R + LoRA	80.28	15.94	9.89
		Router	77.43	10.97	7.68
FiRST (Ours)	R + LoRA	77.5	14.65	10.45	
Router		75.6	9.39	6.92	

Skip (%)	Model Type		BERT	R-1	R-L
0	Original Model	wLoRA	84.89	28.37	17.02
		Base	71.85	19.34	12.00
15	Skip Decode	R + LoRA	83.20	21.71	13.74
		Router	80.97	9.74	6.87
	Random Skip	R + LoRA	79.52	20.18	12.10
		Router	68.20	10.10	7.10
	Unified Skip	R + LoRA	81.53	18.89	11.72
		Router	70.01	12.49	8.68
FiRST (Ours)	R + LoRA	83.17	26.47	16.79	
Router		70.98	16.47	10.51	
24	Skip Decode	R + LoRA	78.55	15.83	6.74
		Router	76.91	13.29	8.86
	Random Skip	R + LoRA	80.00	16.33	10.07
		Router	67.88	8.49	5.96
	Unified Skip	R + LoRA	79.31	15.88	10.69
		Router	68.86	9.17	6.97
FiRST (Ours)	R + LoRA	80.25	21.28	13.89	
Router		69.17	12.36	8.43	
28	Skip Decode	R + LoRA	70.69	8.76	6.74
		Router	40.99	2.05	1.23
	Random Skip	R + LoRA	79.45	14.69	9.23
		Router	67.48	8.14	5.64
	Unified Skip	R + LoRA	78.57	11.74	7.47
		Router	68.23	8.12	5.66
FiRST (Ours)	R + LoRA	77.48	15.98	11.14	
Router		67.14	8.09	6.00	

Table 11: Quality Analysis on Summarization (CNN/DM dataset) on LLaMA-3-8B (left) and LLaMA-3.2-3B (right): BERT F1, Rouge-1 and Rouge-L scores are reported for varying skipping levels. Note that R + LoRA corresponds to Router Augmentation followed by LoRA fine-tuning (in the proposed FiRST framework) and wLoRA stands for Base Model with LoRA fine-tuning. FiRST with fine-tuning, improves upon Unified Skipping for all skipping levels on both Rouge-1 and Rouge-L and is competitive on BERT F1.

Model Type	~ Skipping (%)	Eng→De TPOT	Eng→Zh TPOT
Base + LoRA	0	1x	1x
R + LoRA	15	0.90x	0.88x
R + LoRA	25	0.82x	0.83x
R + LoRA	35	0.69x	0.68x

Model Type	~ Skipping (%)	CNN/DM TPOT
Base + LoRA	0	1x
R + LoRA	15	0.88x
R + LoRA	20	0.81x
R + LoRA	27	0.76x

Table 12: TPOT variation of LLaMA-3-8B on WMT (left) and CNN/DM (right) for FiRST. The reported values are relative to the LoRA fine-tuned base model. Fine-tuning improves TPOT and quality significantly.

Model Type	~ Skipping (%)	Eng→De TPOT	Eng→Zh TPOT
Base + LoRA	0	1x	1x
R + LoRA	15	0.90x	0.91x
R + LoRA	25	0.78x	0.75x
R + LoRA	35	0.69x	0.74x

Model Type	~ Skipping (%)	CNN/DM TPOT
Base + LoRA	0	1x
R + LoRA	15	0.88x
R + LoRA	24	0.79x
R + LoRA	28	0.77x

Table 13: TPOT variation of LLaMA-3.2-3B on WMT (left) and CNN/DM (right) for FiRST. The reported values are relative to the LoRA fine-tuned base model. Fine-tuning improves TPOT and quality significantly.

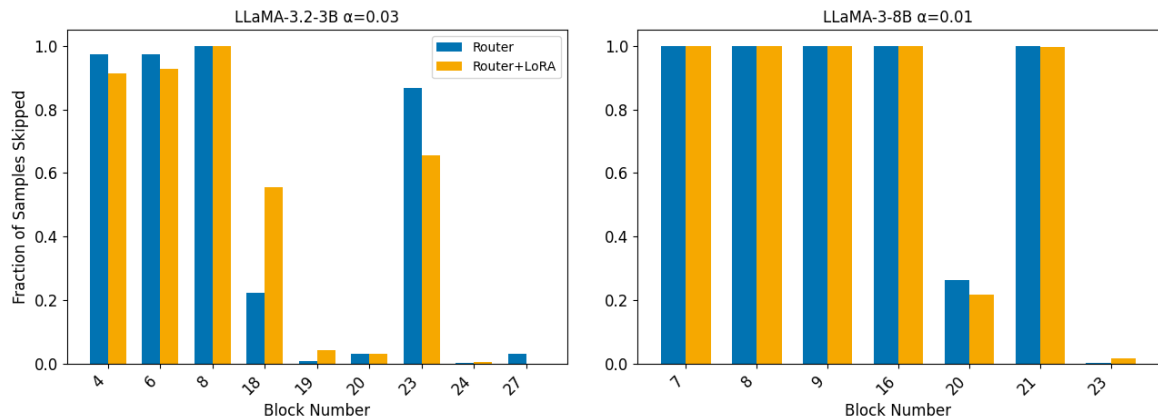


Figure 8: Comparison of LLaMA-3.2-3B (left) and LLaMA-3-8B (right) at 15% skipping rate on English-to-Chinese Machine Translation Task. The graph shows how different layers contribute to the skipping behavior for the same dataset. Layers with no skipping, indicated by a 0% skipping rate, are not represented in the plot.