# Dynamic Vocabulary Pruning in Early-Exit LLMs

**Jort Vincenti**[1,*]    **Karim Abdel Sadek**[1,3,*]    **Joan Velja**[1,*]    **Matteo Nulli**[1,*]
**Metod Jazbec**[1,2]

[1]University of Amsterdam  [2]UvA-Bosch Delta Lab  [3]Krueger AI Safety Lab (KASL)

## Abstract

Increasing the size of large language models (LLMs) has been shown to lead to better performance. However, this comes at the cost of slower and more expensive inference. Early-exiting is a promising approach for improving the efficiency of LLM inference by enabling next token prediction at intermediate layers. Yet, the large vocabulary size in modern LLMs makes the confidence estimation required for exit decisions computationally expensive, diminishing the efficiency gains. To address this, we propose dynamically pruning the vocabulary at test time for each token. Specifically, the vocabulary is pruned at one of the initial layers, and the smaller vocabulary is then used throughout the rest of the forward pass. Our experiments demonstrate that such post-hoc dynamic vocabulary pruning improves the efficiency of confidence estimation in early-exit LLMs while maintaining competitive performance.

## 1 Introduction

Large language models (LLMs) are increasingly being adopted due to their impressive performance and their few-shot ability to adapt to new tasks [3]. However, their growing size results in slow and costly inference. This is particularly limiting in environments with constrained resources or low-latency requirements (e.g., on-device). The push for more efficient LLM implementations is further motivated by growing concerns over their carbon footprint [10]. As a result, making LLMs more efficient at test time has recently received a lot of attention [2, 21, 20, 9]. One promising paradigm for more efficient inference is *early-exiting* [16]. In this case, the forward pass is accelerated by enabling the model to yield a prediction (token) at intermediate layers, rather than passing through all the layers as is traditionally done.

A key component of early-exit models is the *confidence score*, computed at every candidate exit, which determines whether the current prediction is of sufficient quality to terminate the forward pass and return the prediction. While various confidence measures have been proposed, most are derived from the predictive distribution at the given exit (e.g., maximum softmax probability). However, this poses a problem when applying early-exiting to LLMs [4, 14, 1, 17], where obtaining the predictive distribution requires mapping the current hidden representation to the vector of logits over all possible tokens. Given the large vocabulary sizes used in modern LLMs ($\approx$ 30-256K) [19, 15], such confidence estimation introduces significant computational overhead. This is one of the main reasons behind the previously observed paradox, where early-exiting in LLMs resulted in less efficient inference compared to standard, non-accelerated models (both in terms of FLOPs [14] and latency [1]), thereby defeating its original purpose.

In this work, we improve the efficiency of confidence estimation in early-exit LLMs. Specifically, we propose to map the hidden representation of the model to the full vocabulary only at the first couple of initial candidate exits, and use the resulting predictive distribution to identify the top $K$ most likely tokens. We then prune the weight matrix (which maps hidden representations to logits over

---

*Equal contributions. Corresponding author: <m.jazbec@uva.nl>

tokens) based on the most likely tokens found and use the pruned weights at all subsequent candidate exits (Figure 1). Our design is motivated by the empirical observation that the token predicted at the final layer is among the top tokens already in the early layers of the forward pass (Figure 2). In our experiments, we demonstrate that *dynamic vocabulary pruning* improves the FLOPs and time efficiency of confidence estimation in early-exit LLMs while preserving competitive performance. Importantly, our design is lightweight, as it is entirely post-hoc and requires no finetuning or the introduction of new model parameters

## 2  Preliminaries

Let $\mathcal{Y}$ denote the vocabulary (or token) space, with size $|\mathcal{Y}| = d_{\text{vocab}}$. Further, for $x_i \in \mathcal{Y}$, let $(x_1, \ldots, x_t)$ represent the input sequence, comprising both the tokens in the prompt and those generated upon time $t$ by the model.

**Autoregressive Decoding in LLMs**  To predict the next token in the sequence, most modern language models employ the transformer architecture [18]. In a transformer model, the input sequence is passed through $L$ layers, each consisting of a multi-head attention and a feed-forward block, yielding a sequence of hidden representations $\{\mathbf{h}_t^\ell\}_{\ell=1}^L$, $\mathbf{h}_t^\ell \in \mathbb{R}^{d_{\text{model}}}$. After processing through all layers, the final next token distribution is obtained via

$$p\left(x_{t+1}|\mathbf{h}_t^L\right) = \text{softmax}(\mathbf{W}\mathbf{h}_t^L).$$

$\mathbf{W} \in \mathbb{R}^{d_{\text{vocab}} \times d_{\text{model}}}$ is a weight matrix, also referred to as the *unembedding* matrix, that projects the final hidden state $\mathbf{h}_t^L$ back to the token space $\mathcal{Y}$. The newly predicted token $x_{t+1}$ is then added to the input sequence, and the (autoregressive) generation process is repeated until termination.

**Early-Exiting in LLMs**  Observe how decoding in LLMs, as introduced above, requires passing through all L layers for every token in the generated sequence, resulting in a slow inference process. To mitigate this, early-exiting (EE) mechanisms have been proposed [4, 14], allowing the model to predict tokens at intermediate layers if sufficiently confident. Specifically, for each layer $\ell$, a confidence score $c_t^\ell \in [0, 1]$ and an exiting threshold $\lambda_t^\ell \in [0, 1]$ are defined. The early prediction is returned as soon as the confidence at the current layer exceeds the threshold:

$$x_{t+1} := \begin{cases} \arg\max p\left(x_{t+1}|\mathbf{h}_t^1\right) & \text{if } c_t^1 \geq \lambda_t^1, \\ \arg\max p\left(x_{t+1}|\mathbf{h}_t^2\right) & \text{if } c_t^2 \geq \lambda_t^2, \\ \vdots & \vdots \\ \arg\max p\left(x_{t+1}|\mathbf{h}_t^L\right) & \text{otherwise.} \end{cases} \tag{1}$$

Note that it is common to reuse the final weight matrix $\mathbf{W}$ at earlier exits [14, 5], i.e., $p\left(x_{t+1}|\mathbf{h}_t^\ell\right) = \text{softmax}(\mathbf{W}\mathbf{h}_t^\ell), \forall \ell = 1, \ldots, L$, which avoids instantiating a separate unembedding matrix at each exit and prevents introducing a significant number of additional model parameters. Moreover, for simplicity, it is common to assume a fixed and shared threshold $\lambda$ across all exits and tokens [8].

## 3  Dynamic Vocabulary Pruning

**Softmax Based Confidence Measures**  As introduced in Section 2, a confidence measure is necessary to determine whether the model's current prediction is of sufficient quality to terminate the forward pass and return an early prediction. Most commonly, the so-called softmax based measures are used, e.g. the maximum softmax probability $c_t^\ell = \max p(x_{t+1}|\mathbf{h}_t^\ell)$. However, this requires computations involving the full unembedding matrix $\mathbf{W}$ at every exit, which is expensive due to the large $d_{\text{model}}$ and $d_{\text{vocab}}$ used in modern LLMs.[2] While this may be less concerning for latency—since the execution of the next transformer block can proceed in parallel with the confidence estimation—it still reduces the overall efficiency of the forward pass. For example, in CALM [14], the authors report that their early-exit model with softmax confidence is approximately twice as expensive in

---

[2]Confidence measures based directly on the hidden states $\mathbf{h}_t^\ell$ have also been explored, but they have been shown to result in slower exiting compared to softmax-based scores [14].

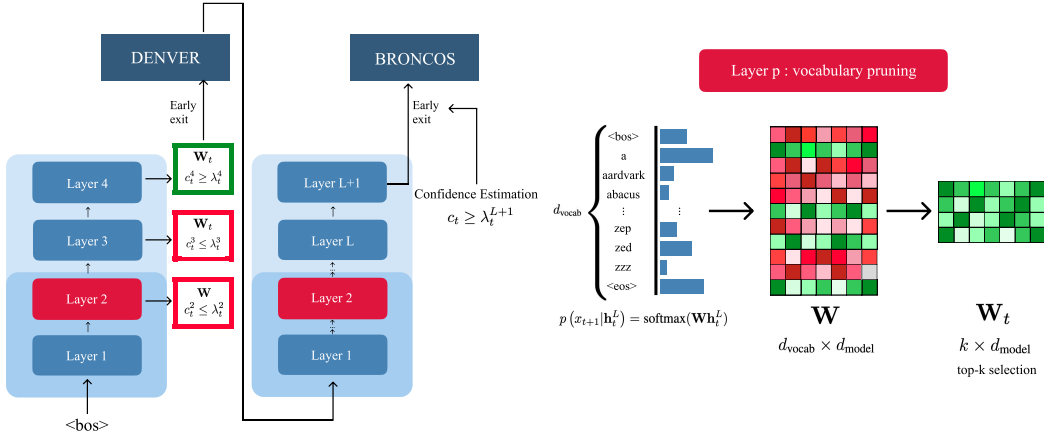**Q**: Which NFL team represented the AFC at Superbowl 50?



Figure 1: *Left*: Illustration of our vocabulary pruning setup in Transformer models during inference. The model evaluates the input question with an Early Exiting objective where the vocabulary is reduced at a fixed layer $p = 2$ in the reference figure. At each layer $\ell$, the model computes a confidence estimation $c_t^\ell$ and compares it against a threshold $\lambda_t^\ell$. When the model achieves sufficient confidence about the token to predict at layer $\ell + 1$, the token is returned. *Right*: Visualization of our proposed pruning mechanism. At exit $p$, we first identify the top $K$ most likely tokens, which are used to subsample the rows of the unembedding matrix $\mathbf{W}$. The resulting pruned matrix $\mathbf{W}_t$ is then used for confidence estimation at all subsequent exits.

terms of FLOPs compared to a static model (i.e., without early exiting), despite requiring around 50% fewer layers per token on average (see Table 2 in [14]). This can make early-exiting impractical, especially in scenarios where FLOPs are a critical constraint (e.g., on device).

**Dynamic Vocabulary Pruning**   To reduce the overhead of confidence estimation in early-exit LLMs, we investigate whether the full computation with $\mathbf{W}$ is indeed necessary at every candidate exit. In particular, we study how quickly the token predicted after passing through all the layers appears among the most likely tokens at earlier layers. As depicted in Figure 2 , we note that this occurs quite early in the forward pass. For example, in the case of the CALM model [14] on the SQuAD dataset, we observe that the token predicted at the last layer appears among the top $10$ most likely tokens already at the 2nd layer in $95\%$ of cases.[3] This suggests that mapping to the full vocabulary becomes redundant after a certain (early) layer.

We make use of this empirical observation in the design of our pruning solution. Specifically, we propose to map the hidden states to the full vocabulary only up to and including exit $p$ (e.g., $p = 1$ or $p = 2$). Then, we use the logits vector $\mathbf{l}_t^p = \mathbf{W}\mathbf{h}_t^p \in \mathbb{R}^{d_{\text{vocab}}}$ to identify the top $K$ most likely tokens and use those to *prune* the embedding matrix $\mathbf{W}$ (by selecting the rows associated with the indices of the most likely tokens, see Figure 1 ). We denote the pruned matrix as $\mathbf{W}_t \in \mathbb{R}^{K \times d_{\text{model}}}$ and use it to compute the confidence at all subsequent layers. The index $t$ in $\mathbf{W}_t$ highlights the dynamic nature of our pruning, i.e., it is performed independently for each token in the generated sequence. Since $K \ll d_{\text{vocab}}$, the cost of confidence estimation is significantly reduced.

To determine the optimal pruning hyperparameters ($p$ and $K$), we suggest using a small calibration dataset and finding the smallest values for which the performance drop remains negligible. We leave the incorporation of more principled selection mechanisms [8] for future work.

---

[3]We find that early-exit finetuning (see Eq. (6) in [14]) is important for ensuring faster token convergence. See Appendix B for more details.

# 4 Experiments

We closely follow the experimental setup of Schuster et al. [14]. Specifically, we use the `T5-large` model [12] and consider the tasks of question answering (SQuAD [13]) and text summarization (SamSum [6]). As a baseline, we use the CALM model [14] with (full) softmax confidence estimation. Our code is publicly available[4] and we provide further implementation details in Appendix A.

The results are presented in Table 1. First, we observe that for both tasks, our proposed Dynamic Vocabulary Pruning (DVP) either matches the baseline or incurs only a negligible performance drop. Under the same conditions, it outperforms the softmax implementation in CALM [14], in terms of FLOPs and time required for exit decisions. For instance, on the SQuAD dataset, using a conservative exit threshold ($\lambda = 0.99$), our DVP achieves the same F1 score (90.6) while requiring $\sim$ 7x fewer FLOPs than the full softmax baseline. Importantly, unlike other FLOP-efficient confidence measures (e.g., hidden state saturation from [14]), our DVP does not require evaluating additional blocks/layers, as evidenced by the similar average exit block indices compared to the baseline. This observation confirms that the pruned vocabulary terms are indeed the ones not usually predicted by the



Figure 2: Rank (log-scale) of the final predicted token across model exits/layers on SQuAD [13] and SamSum [6] using the early-exit version of the `T5-large` model [1]. We observe a clear trend of very early layers showing a low average rank for the final predicted tokens, which motivates our dynamic vocabulary pruning approach.

model. Moreover, while not the primary focus of our work, it is encouraging that DVP also results in reduced latency (i.e., shorter time required to compute exit confidence). Overall, these results suggest that our dynamic vocabulary pruning method effectively addresses the high cost of confidence estimation in early-exit LLMs with little to no impact on overall performance.
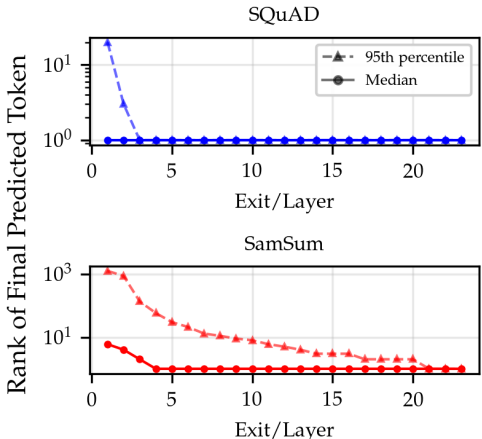
| Dataset | Conf. $\lambda$ | Method | Score ($\uparrow$) | FLOPs/Token ($\downarrow$) | Avg. Exit ($\downarrow$) | Conf. Time (s) ($\downarrow$) |
|---|---|---|---|---|---|---|
| **SQuAD** [13] | 0.6 | CALM | **87.5** | $2.21 \times 10^8$ | <u>2.4</u> | 44.5 |
| | | + DVP (ours) | 87.4 | $\mathbf{1.97 \times 10^8}$ | <u>2.4</u> | **40.8** |
| | $\approx 0.99$ | CALM | <u>90.6</u> | $13.91 \times 10^8$ | 20.9 | 499.9 |
| | | + DVP (ours) | <u>90.6</u> | $\mathbf{1.99 \times 10^8}$ | **20.8** | **413.1** |
| **SamSum** [6] | 0.6 | CALM | **33.8** | $4.21 \times 10^8$ | 5.5 | 90.0 |
| | | + DVP (ours) | 33.7 | $\mathbf{2.01 \times 10^8}$ | **5.4** | **81.0** |
| | $\approx 0.99$ | CALM | <u>43.1</u> | $11.13 \times 10^8$ | 16.5 | 162.0 |
| | | + DVP (ours) | <u>43.1</u> | $\mathbf{2.12 \times 10^8}$ | **16.4** | **136.0** |

Table 1: Summary of efficiency gains for our dynamic vocabulary pruning (DVP) compared to CALM [14] for two different exiting thresholds $\lambda$ (0.6 and $\approx 0.99$). To measure the performance quality, we report F1 score for SQuAD [13] and Rouge-L metric for SamSum [6]. Additionally, we outline the amount of FLOPs per generated token and average early-exit layer across generated tokens (note that the full `T5-large` model has 24 layers). We also report the total time spent on confidence estimation for the entire test set.

# 5 Conclusion & Future Work

Our work tackles the high cost of confidence estimation in early-exit LLMs, which arises from large vocabulary sizes. By dynamically pruning the vocabulary for every generated token, we demonstrate

---

[4]https://github.com/MatteoNulli/Vocabulary_pruning/tree/main

that efficient confidence computation is achievable without compromising performance. Our proposed vocabulary pruning is completely post-hoc, making it nicely compatible with existing pretrained early-exit LLMs. We hope our findings encourage a reconsideration of the trend towards sacrificing model adaptivity (i.e., reducing the number of possible exits [1]) due to the growing computational cost of exiting decisions. In future work, it would be valuable to validate our approach on other early-exit LLMs [17] and explore more advanced pruning mechanisms (e.g., using product-of-experts ensembles across exits [7]) beyond the simple top-K strategy used here. Future work could also investigate the impact of dynamic vocabulary pruning on confidence calibration [11].

## References

[1] S. Bae, J. Ko, H. Song, and S.-Y. Yun. Fast and robust early-exiting framework for autoregressive language models with synchronized parallel decoding, 2023.

[2] G. Bai, Z. Chai, C. Ling, S. Wang, J. Lu, N. Zhang, T. Shi, Z. Yu, M. Zhu, Y. Zhang, et al. Beyond efficiency: A systematic survey of resource-efficient large language models. *arXiv preprint arXiv:2401.00625*, 2024.

[3] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. Language models are few-shot learners, 2020.

[4] M. Elbayad, J. Gu, E. Grave, and M. Auli. Depth-adaptive transformer. *arXiv preprint arXiv:1910.10073*, 2019.

[5] M. Elhoushi, A. Shrivastava, D. Liskovich, B. Hosmer, B. Wasti, L. Lai, A. Mahmoud, B. Acun, S. Agarwal, A. Roman, et al. Layer skip: Enabling early exit inference and self-speculative decoding. *arXiv preprint arXiv:2404.16710*, 2024.

[6] B. Gliwa, I. Mochol, M. Biesek, and A. Wawer. Samsum corpus: A human-annotated dialogue dataset for abstractive summarization. *arXiv preprint arXiv:1911.12237*, 2019.

[7] M. Jazbec, J. Allingham, D. Zhang, and E. Nalisnick. Towards anytime classification in early-exit architectures by enforcing conditional monotonicity. *Advances in Neural Information Processing Systems*, 36, 2024.

[8] M. Jazbec, A. Timans, T. H. Veljković, K. Sakmann, D. Zhang, C. A. Naesseth, and E. Nalisnick. Fast yet safe: Early-exiting with risk control. *arXiv preprint arXiv:2405.20915*, 2024.

[9] S. Kim, C. Hooper, T. Wattanawong, M. Kang, R. Yan, H. Genc, G. Dinh, Q. Huang, K. Keutzer, M. W. Mahoney, et al. Full stack optimization of transformer inference: a survey. *arXiv preprint arXiv:2302.14017*, 2023.

[10] L. Lannelongue, J. Grealey, and M. Inouye. Green algorithms: quantifying the carbon footprint of computation. *Advanced science*, 8(12):2100707, 2021.

[11] L. Meronen, M. Trapp, A. Pilzer, L. Yang, and A. Solin. Fixing overconfidence in dynamic neural networks. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2680–2690, 2024.

[12] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. URL http://jmlr.org/papers/v21/20-074.html.

[13] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.

[14] T. Schuster, A. Fisch, J. Gupta, M. Dehghani, D. Bahri, V. Q. Tran, Y. Tay, and D. Metzler. Confident adaptive language modeling, 2022.

[15] C. Tao, Q. Liu, L. Dou, N. Muennighoff, Z. Wan, P. Luo, M. Lin, and N. Wong. Scaling laws with vocabulary: Larger models deserve larger vocabularies. *arXiv preprint arXiv:2407.13623*, 2024.

[16] S. Teerapittayanon, B. McDanel, and H.-T. Kung. Branchynet: Fast inference via early exiting from deep neural networks. In *2016 23rd international conference on pattern recognition (ICPR)*, pages 2464–2469. IEEE, 2016.

[17] N. Varshney, A. Chatterjee, M. Parmar, and C. Baral. Accelerating llm inference by enabling intermediate layer decoding. *arXiv preprint arXiv:2310.18581*, 2023.

[18] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[19] P. Villalobos, A. Ho, J. Sevilla, T. Besiroglu, L. Heim, and M. Hobbhahn. Will we run out of data? limits of llm scaling based on human-generated data, 2024.

[20] M. Xu, W. Yin, D. Cai, R. Yi, D. Xu, Q. Wang, B. Wu, Y. Zhao, C. Yang, S. Wang, et al. A survey of resource-efficient llm and multimodal foundation models. *arXiv preprint arXiv:2401.08092*, 2024.

[21] Z. Zhou, X. Ning, K. Hong, T. Fu, J. Xu, S. Li, Y. Lou, L. Wang, Z. Yuan, X. Li, et al. A survey on efficient inference for large language models. *arXiv preprint arXiv:2404.14294*, 2024.

# Appendix

## A Implementation Details

We ran our experiments on 1x Nvidia A100 80GB - SMX4 GPU. Our code is available at `https://github.com/MatteoNulli/Vocabulary_pruning/tree/main`.

We report all the relevant early-exiting hyperparameters for our experiments in Table 2. Our DVP approach introduces $p$ and $K$ which represent the pruning exit index and the pruned vocabulary size, respectively. The `top-2 diff` strategy indicates that the exit confidence $c_t^\ell$ is computed as the difference between the probabilities of the top two tokens. The `decaying` threshold $\lambda_t$ means that the exit threshold decreases for later tokens in the generated response (see Eq. (5) in [14]).

|  | SQuAD | | SamSum | |
|---|---|---|---|---|
|  | CALM | DVP | CALM | DVP |
| $p$ | - | 2 | - | 2 |
| $K$ | - | 64 | - | 512 |
| $c_t^\ell$ | `top-2 diff` | `top-2 diff` | `top-2 diff` | `top-2 diff` |
| $\lambda_t$ | `static` | `static` | `decaying` $(\tau = 4)$ | `decaying` $(\tau = 4)$ |

Table 2: Main early-exit hyperparameters used in our experiments.

## B Additional Experiments

In Section 3, we reported that, for an early-exit LLM like CALM [14], the token predicted at the final layer is often among the top-K predicted tokens quite early in the process. Here, we investigate the effect of adapting the unembedding matrix $\mathbf{W}$ to intermediate representations $\mathbf{h}_t^\ell$ through early-exit finetuning (see Eq. (6) in [14]). The results, displayed in Figure 3, show that the T5 model [12] without early-exit finetuning exhibits slower convergence compared to the T5 model that has undergone early-exit finetuning (which corresponds to the CALM model). This finding is important for our dynamic vocabulary pruning proposal, as faster convergence enables the selection of lower values for pruning parameters ($p$ and $K$), resulting in larger efficiency savings.
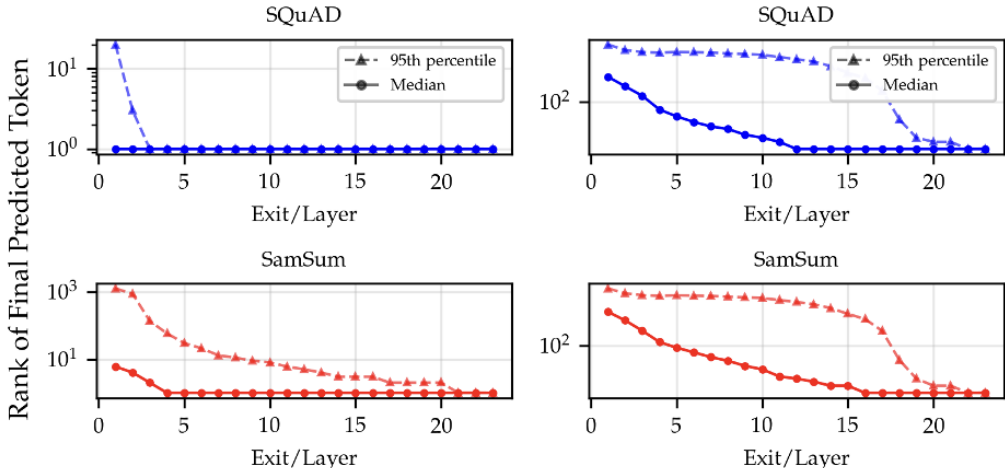


Figure 3: Rank (log-scale) of the final predicted token across model exits/layers on SQuAD [13] and SamSum [6]. *Left:* Results based on CALM [14], the early-exit version of the `T5-large` model [1]. These are the same results as those shown in Figure 2, included here for easier comparison. *Right:* Results based on the `T5-large` model [12], where the non-adapted original unembedding matrix is used at intermediate layers to facilitate early-exiting.