

# SPARK TRANSFORMER: HOW MANY FLOPS IS A TOKEN WORTH?

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

This work introduces *Spark Transformer*, an architectural variant of the Transformer model that drastically reduces the FLOPs count while maintaining comparable quality and an identical parameter count. This reduction is achieved by introducing sparse activations in both the feedforward network (FFN) and the Attention mechanism. In the FFN, this sparsity engages only a subset of parameters for each input. In the Attention mechanism, it limits the number of tokens that each token attends to. We achieve this sparsity through *statistical top-k*, a lightweight approximate algorithm that is well-suited for accelerator hardware and minimizes training slowdown. Furthermore, Spark Transformer incorporates dedicated predictors to identify the activated entries. These predictors are formed by allocating a portion of the model’s parameters and are trained jointly with the rest of the model. This approach distinguishes Spark Transformer from existing methods that introduce sparsity and predictors post-training, which often leads to increased training costs, additional model parameters, and complex modifications to the model architecture. Our Spark Transformer, pretrained using the Gemma 2 recipe, achieves competitive performance on standard benchmarks while exhibiting significant sparsity. Specifically, it utilizes only 8% nonzeros in the FFN activation and attends to a maximum of 256 tokens. This results in a  $3.1\times$  reduction in FLOPs, yielding a  $1.70\times$  speedup for prefill and a  $1.79\times$  speedup for decoding on a 16-core CPU VM.

## 1 INTRODUCTION

The machine learning landscape has witnessed a surge in large-scale Transformer models (Anil et al., 2023; Almazrouei et al., 2023; Dubey et al., 2024; Adler et al., 2024), pushing the boundaries of language understanding and generation. However, the pursuit of scale is often constrained not by limitations in model quality, but by the escalating computational costs (Sharir et al., 2020; Patterson et al., 2021) associated with increasing parameter counts (Kaplan et al., 2020). This challenge is further exacerbated by the development of models that handle increasingly long context inputs (Reid et al., 2024), where computational demands grow proportionally with context length.

*Sparse activation* is a popular approach for addressing the computational challenges posed by both large model size and long context length. To handle large models, sparse activation reduces costs by engaging only a small subset of model parameters for each input. This approach has gained significant interest following the discovery of an intriguing phenomenon: the feed-forward networks (FFNs) in classic Transformers like T5 (Raffel et al., 2020) and ViT (Dosovitskiy, 2020) exhibit *natural* activation sparsity (Zhang et al., 2022; Li et al., 2022). In other words, these models demonstrate sparsity without explicit enforcement. This inherent sparsity, reaching remarkable levels like 3% nonzeros in T5, presents an opportunity for substantial efficiency gains with minimal modifications to the model architecture.

Unfortunately, this natural sparsity is absent in the latest generation of models (Jiang et al., 2023; Gemma Team, 2024; Dubey et al., 2024), which have adopted *gated non-ReLU* activation functions (Dauphin et al., 2017). To re-introduce sparsity, recent work has explored extra pretraining steps, either by switching back to ReLU activations (Mirzadeh et al., 2023; Zhang et al., 2024) or by adding top- $k$  thresholding (Yerram et al., 2024; Song et al., 2024a;b). This is often combined with training a low-cost predictor to identify activated parameters, a crucial step for maximizing efficiency gains (Liu et al., 2023; Zeng et al., 2023; Song et al., 2023; Yerram et al., 2024). However,

Table 1: **FLOPs per token comparison: Spark Transformer vs. standard Transformer.** In a standard Transformer with model dimension  $d_{\text{model}}$ , we assume multi-head attention where the sum of head dimensions equals  $d_{\text{model}}$ , and an FFN with non-gated activation and width  $d_{\text{ff}}$ . Here,  $n_{\text{ctx}}$  represents the context length for the target token. The computational cost is primarily determined by the FFN (assuming  $d_{\text{ff}} \gg d_{\text{model}}$ , which is typical) and the attention dot product (assuming a long context length). Spark Transformers introduce sparsity parameters,  $k_{\text{ff}}$  and  $k_{\text{attn}}$ , to reduce FLOPs. Setting  $k_{\text{ff}} = 8\% \times d_{\text{ff}}$  and  $k_{\text{attn}} = 256$  achieves a  $3.2\times$  FLOPs reduction in the FFN, a  $4\times$  reduction in the attention dot product, and a  $3.1\times$  reduction overall (assuming  $n_{\text{ctx}} = 8k$ ) for Gemma-2B.

| Operation                   | FLOPs per Token <sup>1</sup>      |  |
|-----------------------------|-----------------------------------|--|
|                             | Standard Transformer              | Spark Transformer (Ours)   |
| FFN                         | $2d_{\text{model}}d_{\text{ff}}$  | $0.5d_{\text{model}}d_{\text{ff}} + 1.5d_{\text{model}}k_{\text{ff}}$                            |
| Attention dot product       | $2d_{\text{model}}n_{\text{ctx}}$ | $0.5d_{\text{model}}n_{\text{ctx}} + 1.5d_{\text{model}}\min\{k_{\text{attn}}, n_{\text{ctx}}\}$ |
| Attention linear projection | $4d_{\text{model}}^2$             | $4d_{\text{model}}^2$  |

these approaches not only complicate the training procedure and incur extra training costs, but also introduce additional parameters (for the predictor) and have yet to demonstrate high sparsity levels without compromising model quality.

Sparse activation for attention, often called sparse attention, faces a similar challenge. Sparsity is used to efficiently handle long context input by limiting the number of tokens each token attends to. A straightforward approach is top- $k$  attention (Gupta et al., 2021), which applies a top- $k$  mask to the attention coefficients. This can be combined with a low-cost predictor to maximize efficiency (Ribar et al., 2023; Yang et al., 2024; Lee et al., 2024c). However, achieving high sparsity and a predictor without complicated procedure and sacrificing quality remains a challenge.

**Contributions.** This work introduces *Spark Transformer*, an architectural variant of Transformer that achieves both high activation sparsity and low-cost prediction in both the FFN and Attention mechanisms. Notably, Spark Transformer can be trained in a single stage without requiring separate post-processing and maintains quality without introducing additional parameters. This makes it a suitable drop-in replacement for standard Transformer models. Because the FFN and attention components dominate the computational cost in large Transformers with long contexts, Spark Transformer drastically reduces the overall FLOP count for decoding a token (see Table 1).

We pretrain a Spark Transformer using the Gemma-2 recipe (Gemma Team, 2024), resulting in a model we call Spark Gemma-2. Evaluation on standard benchmarks demonstrates that Spark Gemma-2 closely matches the quality of Gemma-2, even with a high degree of sparsity: only 8% activated entries in the FFN and a maximum of 256 attended tokens in attention (see Table 2). This sparsity leads to a  $3.1\times$  reduction in overall FLOPs compared to Gemma-2. Using this sparsity, we evaluate model efficiency with `gemma.cpp` (Google Gemma.cpp, 2024), a C++ inference engine optimized for serving Gemma models on CPUs, and observe a speedup of up to  $1.79\times$  (see Figure 3). Notably, on a readily available 4-core cloud VM, Spark Gemma-2 achieves a decoding speed of 86 ms per token, surpassing typical human reading speed (Brysbaert, 2019). This increased efficiency enables wider access to high-quality models for users with limited access to high-FLOP devices, such as GPUs and TPUs.

Spark Transformer leverages the interpretation of both FFN and attention as key-value lookup tables (Geva et al., 2021) to provide a unified solution for achieving sparsity and prediction. Specifically, the predictor is obtained by repurposing a subset of the dimensions of the query and key vectors to produce an importance score for each key-value pair (see Section 3). Top- $k$  thresholding is applied to these scores to identify the activated keys. In particular, standard top- $k$  thresholding requires performing a sorting, which is inefficient particularly on training accelerators. To address this, we introduce *statistical top- $k$* , a linear complexity algorithm for approximate nearest neighbor search (see Section 2) based on fitting a Gaussian distribution to the activation entries and estimating a threshold that yields the top entries. **While ideas similar to statistical top- $k$  have been used (Shi et al., 2019; M Abdelmoniem et al., 2021) for the problem of distributed training (Lin et al., 2018), we are the first to introduce, adapt, and verify its effectiveness for activation sparsity.**

<sup>1</sup>Please refer to Section 3.1 and Section 3.2 for the calculation of FLOPs for FFN and Attention, respectively. We omit non-leading-order terms (e.g., those arising from embedding, normalization, and nonlinear layers) and exclude the number of layers as a common multiplier.

## 2 STATISTICAL TOP-K

This section introduces Statistical-Top<sub>k</sub>, an approximate algorithm for obtaining the  $k$  largest entries of an input vector. Recall that the *soft-thresholding operator* is defined for an arbitrary vector  $\mathbf{x} \in \mathbb{R}^d$  and a scalar threshold  $\theta \in \mathbb{R}$  as

$$\text{Soft-Threshold}(\mathbf{x}, \theta) \stackrel{\text{def}}{=} \max\{\mathbf{x} - \theta \cdot \mathbf{1}, \mathbf{0}\} \in \mathbb{R}^d, \quad (1)$$

where  $\mathbf{1}$  and  $\mathbf{0}$  are  $d$ -dimensional vectors with all entries equal to 1 and 0, respectively. The soft-thresholding operator shifts each entry of  $\mathbf{x}$  to the left by  $\theta$  and then thresholds the result at zero.

We define Statistical-Top<sub>k</sub> as the following mapping from  $\mathbb{R}^d$  to  $\mathbb{R}^d$ :

$$\text{Statistical-Top}_k(\mathbf{x}) \stackrel{\text{def}}{=} \text{Soft-Threshold}(\mathbf{x}, \theta(\mathbf{x}, k)), \text{ where } \theta(\mathbf{x}, k) \stackrel{\text{def}}{=} \text{mean}(\mathbf{x}) + \text{std}(\mathbf{x}) \cdot Q(1 - \frac{k}{d}) \quad (2)$$

Here,  $\text{mean}(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{d} \sum_{i=1}^d x_i$  and  $\text{std}(\mathbf{x}) \stackrel{\text{def}}{=} \sqrt{\frac{1}{d-1} \sum_{i=1}^d (x_i - \text{mean}(\mathbf{x}))^2}$  compute the sample mean and standard deviation of the entries of  $\mathbf{x}$ , respectively, and  $Q(\cdot)$  is the quantile function (i.e., inverse of the cumulative distribution function) of the standard Gaussian distribution.

Statistical-Top<sub>k</sub> in Eq. (2) operates by first computing a threshold  $\theta(\mathbf{x}, k)$  such that approximately  $k$  entries of  $\mathbf{x}$  exceed it, and then applying the soft-thresholding operator with this threshold to  $\mathbf{x}$  to obtain a sparse output. We discuss these two components in the next two subsections.

### 2.1 THRESHOLD ESTIMATION

The threshold  $\theta(\mathbf{x}, k)$  in Eq. (2) is designed such that, if the entries of  $\mathbf{x}$  are drawn from a Gaussian distribution, approximately  $k$  out of the  $d$  entries will exceed this threshold. To understand this, let  $\mu$  and  $\sigma$  denote the mean and standard deviation of the underlying Gaussian distribution. Its quantile function is given by  $\mu + \sigma \cdot Q(p)$  for  $p \in (0, 1)$ . Consequently, due to the properties of quantile functions, we expect roughly  $p \cdot d$  entries of  $\mathbf{x}$  to exceed  $\mu + \sigma \cdot Q(1 - p)$ . In practice, since  $\mu$  and  $\sigma$  are unknown, they are replaced with the sample mean  $\text{mean}(\mathbf{x})$  and the sample standard deviation  $\text{std}(\mathbf{x})$ , respectively.

The following theorem formalizes this argument.

**Theorem 1.** *Let  $\mathbf{x} \in \mathbb{R}^d$  be a vector with entries drawn i.i.d. from  $\mathcal{N}(\mu, \sigma^2)$ . For any  $1 \leq k \leq d-1$ , let  $\theta(\mathbf{x}, k)$  be a scalar defined in Eq. (2). Take any  $\delta \in (0, 1)$  and assume  $d \geq \max\{2, \log \frac{6}{\delta}\}$ . With a probability of at least  $1 - \delta$ , the number of entries of  $\mathbf{x}$  that are greater than  $\theta(\mathbf{x}, k)$ , i.e.,  $\text{card}(\{i \in [d] \mid x_i > \theta(\mathbf{x}, k)\})$ , satisfies*

$$\left| \frac{\text{card}(\{i \in [d] \mid x_i > \theta(\mathbf{x}, k)\}) - k}{d} \right| \leq 4 \sqrt{\frac{\log \frac{6}{\delta}}{d}} \left( 1 + \sqrt{-2 \log \min\left\{\frac{k}{d}, 1 - \frac{k}{d}\right\}} \right).$$

Theorem 1 provides a relative error bound between  $k$  and the true number of entries of  $\mathbf{x}$  that exceed  $k$ . This bound is maximized when  $k = 1$  or  $k = d - 1$ . Consequently, the worst-case bound is  $O\left(\sqrt{\frac{\log d \cdot \log \frac{1}{\delta}}{d}}\right)$  which vanishes as  $d$  increases. Notably, the error bound becomes  $O\left(\sqrt{\frac{\log \frac{1}{\delta}}{d}}\right)$  when  $k = \Theta(d)$ , demonstrating even faster convergence.

**Computation cost.** The computation of the threshold  $\theta(\mathbf{x}, k)$  is highly efficient, requiring only  $2d$  FLOPs to compute the mean and standard deviation of the samples. This contrasts sharply with a naive sorting-based approach, which has  $O(d \log d)$  complexity.

While the Gaussian quantile function  $Q(\cdot)$  lacks a closed-form solution, high-precision piecewise approximation algorithms with constant complexity are available in standard software packages like SciPy (Virtanen et al., 2020), readily applicable to our needs.

### 2.2 SPARSIFICATION

Given the threshold  $\theta(\mathbf{x}, k)$ , a straightforward approach to obtain a sparse vector is to set all entries of  $\mathbf{x}$  below the threshold to zero, preserving the remaining values. This operator, sometimes referred to as *hard thresholding* (Blumensath & Davies, 2008), suffers from discontinuity, potentially hindering its suitability for gradient-descent-based training.

To address this, Statistical-Top<sub>k</sub> employs the soft-thresholding operator defined in Eq. (1) (Beck & Teboulle, 2009). This operator first shrinks all entries of  $\mathbf{x}$  by the threshold  $\theta(\mathbf{x}, k)$  and then sets all entries below 0 to 0. Soft thresholding offers the advantages of being continuous and differentiable almost everywhere (except when entries of  $\mathbf{x}$  coincide with  $\theta(\mathbf{x}, k)$ ).

For complete differentiability, one can utilize a smoothing function like the Huber loss (Huber, 1992), defined element-wise on an input  $\mathbf{x}$  as:

$$\text{Huber}(x; \delta) \stackrel{\text{def}}{=} \begin{cases} \frac{1}{2}x^2 & \text{for } |x| < \delta, \\ \delta \cdot (|x| - \frac{1}{2}\delta) & \text{otherwise.} \end{cases} \quad (3)$$

The following theorem establishes the continuous differentiability of the mapping  $\mathbf{x} \mapsto \text{Huber}(\text{Statistical-Top}_k(\mathbf{x}); \delta)/\delta$ :

**Theorem 2.** For any  $\delta > 0$ , the function  $\mathbb{R}^d \rightarrow \mathbb{R}^d$  defined as

$$\text{Huber}(\text{Statistical-Top}_k(\mathbf{x}); \delta) / \delta \quad (4)$$

is continuously differentiable.

Note that Eq. (4) converges to Statistical-Top<sub>k</sub>( $\mathbf{x}$ ) as  $\delta \rightarrow 0$ , since  $\text{Huber}(x; \delta)/\delta \rightarrow |x|$  and Statistical-Top<sub>k</sub>( $\mathbf{x}$ ) is always non-negative. In practice, however, we find that using a non-zero  $\delta$  does not improve model quality, and therefore we set  $\delta = 0$  for simplicity.

Finally, soft thresholding admits a variational form (see, e.g., Parikh et al. (2014)):

$$\text{Soft-Threshold}(\mathbf{x}, \theta) = \arg \min_{\mathbf{z} \geq \mathbf{0}} \theta \|\mathbf{z}\|_1 + \frac{1}{2} \|\mathbf{x} - \mathbf{z}\|_2^2. \quad (5)$$

This formulation seeks a vector  $\mathbf{z}$  that minimizes both its squared  $\ell_2$  distance to the input  $\mathbf{x}$  and its  $\ell_1$  norm, with the threshold  $\theta$  balancing these terms. Given the sparsity-promoting nature of the  $\ell_1$  norm, soft thresholding effectively finds a sparse approximation of the input  $\mathbf{x}$ .

### 2.3 COMPARISON WITH RELATED TOP-*k* OPERATORS

The variational form in Eq. (5) also reveals connections of Statistical-Top<sub>k</sub> with other top-*k* algorithms in the literature. Specifically, Lei et al. (2023) defines a *soft top-k* as

$$\arg \min_{\mathbf{z}} -\theta \cdot H(\mathbf{z}) - \langle \mathbf{z}, \mathbf{x} \rangle, \quad \text{s.t. } \mathbf{z}^\top \mathbf{1} = k, \quad \mathbf{0} \leq \mathbf{z} \leq \mathbf{1}, \quad (6)$$

where  $H(\mathbf{z})$  is the entropy function. Another work (Lou et al., 2024) defines the *SparseK* operator

$$\arg \min_{\mathbf{z}} -H^G(\mathbf{z}) - \langle \mathbf{z}, \mathbf{x} \rangle, \quad \text{s.t. } \mathbf{z}^\top \mathbf{1} = k, \quad \mathbf{0} \leq \mathbf{z} \leq \mathbf{1}, \quad (7)$$

where  $H^G(\mathbf{z})$  is the generalized Gini entropy.

Statistical-Top<sub>k</sub> in the form of Eq. (5), as well as Eq. (6) and Eq. (7), can all be interpreted as finding an output that is *close* to the input subject to a sparsifying regularization. Their major difference lies in the choice of the sparse regularization. That is, Statistical-Top<sub>k</sub> uses  $\ell_1$ , whereas soft top-*k* and SparseK uses entropy and Gini entropy, respectively. The choice of  $\ell_1$  makes Statistical-Top<sub>k</sub> superior in that it has a closed form solution provided by soft-thresholding, which only requires  $d$  FLOPs. In contrast, soft top-*k* and SparseK both do not have closed form solutions and require an iterative algorithm with a FLOP count dependent on the number of iterations. In addition, there is no guarantee that soft top-*k* and SparseK can obtain (approximately)  $k$  nonzero entries as output.

## 3 SPARK TRANSFORMER

This section describes Spark FFN and Spark Attention, the two components of Spark Transformer.

### 3.1 SPARK FFN

FFNs in a standard Transformer are two-layer multi-layer perceptrons that map an input token  $\mathbf{q} \in \mathbb{R}^{d_{\text{model}}}$  to an output

$$\text{FFN}(\mathbf{q}; \mathbf{K}, \mathbf{V}) \stackrel{\text{def}}{=} \mathbf{V} \cdot \sigma(\mathbf{K}^\top \mathbf{q}) \in \mathbb{R}^{d_{\text{model}}}. \quad (8)$$

In above,  $\{\mathbf{K}, \mathbf{V}\} \subseteq \mathbb{R}^{d_{\text{model}} \times d_{\text{ff}}}$  are trainable model parameters, and  $\sigma(\cdot)$  is a nonlinear activation function. We ignore the dependency on layer index to simplify the notations.

Each of the matrix multiplication in Eq. (8) has  $d_{\text{model}} \cdot d_{\text{ff}}$  FLOPs hence overall the computation cost is  $2d_{\text{model}} \cdot d_{\text{ff}}$ . However, previous work shows that when  $\sigma(\cdot)$  is ReLU, the activation map  $\sigma(\mathbf{K}^\top \mathbf{q})$  is very sparse after model training. The sparsity can be used trivially to reduce the computation costs in the calculation of its product with the second layer weight matrix  $\mathbf{V}$  (Li et al., 2022), reducing the overall FLOPs count of FFN to  $d_{\text{model}} \cdot (d_{\text{ff}} + k)$ , where  $k \ll d_{\text{ff}}$  is the number of nonzero entries in the activation. Note that the sparsity cannot be used to reduce the computation costs associated with  $\mathbf{K}$ , which constitute half of the total FLOPs in FFN.

In order to reduce FLOPs count in the first layer of FFN as well, we introduce Spark FFN as follows:

$$\text{Spark-FFN}(\mathbf{q}; \mathbf{K}, \mathbf{V}, k, r) \stackrel{\text{def}}{=} \mathbf{V} \cdot (\sigma(\text{Statistical-Top}_k(\mathbf{K}^\top \mathbf{P}\mathbf{q})) \odot (\mathbf{K}^\top (\mathbf{I} - \mathbf{P})\mathbf{q})) \in \mathbb{R}^{d_{\text{model}}}. \quad (9)$$

In above,  $\{\mathbf{K}, \mathbf{V}\} \subseteq \mathbb{R}^{d_{\text{model}} \times d_{\text{ff}}}$  are trainable parameters as in standard FFNs, and the activation  $\sigma(\cdot)$  is taken to be GELU (Hendrycks & Gimpel, 2016) following Gemma. The Statistical-Top $_k$ , defined in Eq. (2), is introduced for obtaining sparsity, with  $k$  being a hyper-parameter specifying the sparsity level. Finally,  $\mathbf{P}$  is a fixed matrix  $\mathbf{P} \stackrel{\text{def}}{=} \mathbf{1}_r \oplus \mathbf{0}_{d_{\text{model}}-r} \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$  where  $\oplus$  is the direct sum operator and  $r$  is a hyper-parameter. It is introduced so that the term  $\mathbf{K}^\top \mathbf{P}\mathbf{q}$  serves as a low-rank predictor of the location of the nonzero entries, which allows us to obtain efficiency benefits in computing  $\mathbf{K}^\top (\mathbf{I} - \mathbf{P})\mathbf{q}$  and the multiplication with  $\mathbf{V}$ . This is discussed in detail below.

**FLOPs per Token.** Naive implementation of the Spark-FFN has the same number of FLOPs as the vanilla FFN in Eq. (8), i.e.,

$$r \cdot d_{\text{ff}} + (d_{\text{model}} - r) \cdot d_{\text{ff}} + d_{\text{model}} \cdot d_{\text{ff}} = 2d_{\text{model}} \times d_{\text{ff}} \quad (10)$$

where the three terms are from  $\mathbf{K}^\top \mathbf{P}\mathbf{q}$ ,  $\mathbf{K}^\top (\mathbf{I} - \mathbf{P})\mathbf{q}$ , and the multiplication with  $\mathbf{V}$ , respectively. In Spark-FFN, one may first compute the term  $\mathbf{K}^\top \mathbf{P}\mathbf{q}$  as a low-rank predictor. After passing its output through Statistical-Top $_k$ , which selects approximately the  $k$  most important entries, followed by the activation function  $\sigma(\cdot)$ , we obtain a sparse output. Importantly, after obtaining the sparse output there is no need to perform the full computation of the other two matrix multiplications in Eq. (9), i.e.,  $\mathbf{K}^\top (\mathbf{I} - \mathbf{P})\mathbf{q}$  and the multiplication with  $\mathbf{V}$ . Instead, one can perform a sparse matrix multiplication with a drastically reduced FLOPs count:

$$r \cdot d_{\text{ff}} + (d_{\text{model}} - r) \cdot k + d_{\text{model}} \cdot k = (d_{\text{ff}} - k) \cdot r + 2d_{\text{model}} \cdot k, \quad (11)$$

which is an increasing function of  $r$ . In other words,  $r$  controls the computation cost. We provide ablation study in the section to show that the best model quality is obtained when  $r \approx \frac{d_{\text{model}}}{2}$ . In this case, the total FLOP count of Spark FFN is approximately  $0.5 \cdot d_{\text{model}} \cdot d_{\text{ff}} + 1.5 \cdot d_{\text{model}} \cdot k$ , which is a 4-times reduction from Eq. (10) when  $k$  is very small.

**Relation to gated activation.** Many of the most recent Transformers, including Gemma 2, use a variant of the standard FFN in Eq. (8) where the activation function is replaced with a gated one:

$$\text{Gated-FFN}(\mathbf{q}; \mathbf{K}_1, \mathbf{K}_2, \mathbf{V}) = \mathbf{V} \cdot (\sigma(\mathbf{K}_1^\top \mathbf{q}) \odot (\mathbf{K}_2^\top \mathbf{q})). \quad (12)$$

Note that when compared with the FFN in Eq. (8) for quality studies,  $d_{\text{ff}}'$  is usually taken to be  $2/3 \cdot d_{\text{ff}}$  to be iso-parameter count (Shazeer, 2020).

Our Spark FFN in Eq. (9) bears some resemblance to Gated FFN in that both have two linear maps in the first layer and one in the second layer. The difference lies in that 1) Spark FFN adds a statistical top- $k$  to obtain sparsity, and 2) the input to the first layers of Spark FFN are obtained from splitting the dimensions of the input. The latter change has the benefit that it offers a convenient means of controlling the number of FLOPs from tuning the choice of  $r$  (see Eq. (11)).

### 3.2 SPARK ATTENTION

In a standard multi-head attention layer, an input  $\mathbf{x} \in \mathbb{R}^{d_{\text{model}}}$  is mapped to a query, a key, and a value vector of dimension  $d_{\text{attn}}$  as  $\mathbf{q}^{(i)} = \mathbf{W}_Q^{(i)} \mathbf{x} \in \mathbb{R}^{d_{\text{attn}}}$ ,  $\mathbf{k}^{(i)} = \mathbf{W}_K^{(i)} \mathbf{x} \in \mathbb{R}^{d_{\text{attn}}}$ ,  $\mathbf{v}^{(i)} = \mathbf{W}_V^{(i)} \mathbf{x} \in \mathbb{R}^{d_{\text{attn}}}$  for each head  $i$ . Here,  $\{\mathbf{W}_Q^{(i)}, \mathbf{W}_K^{(i)}, \mathbf{W}_V^{(i)}\} \subseteq \mathbb{R}^{d_{\text{attn}} \times d_{\text{model}}}$  are trainable weights.

Collecting all the key and value vectors in the context of  $\mathbf{x}$  into  $\mathbf{K}^{(i)} = [\mathbf{k}_1^{(i)}, \dots, \mathbf{k}_{n_{\text{ctx}}}^{(i)}] \in \mathbb{R}^{d_{\text{attn}} \times n_{\text{ctx}}}$  and  $\mathbf{V}^{(i)} = [\mathbf{v}_1^{(i)}, \dots, \mathbf{v}_{n_{\text{ctx}}}^{(i)}] \in \mathbb{R}^{d_{\text{attn}} \times n_{\text{ctx}}}$ , attention conducts the following computation:

$$\text{Attention}(\mathbf{q}; \mathbf{K}, \mathbf{V}) \stackrel{\text{def}}{=} \mathbf{V} \cdot \text{softmax}(\mathbf{K}^\top \mathbf{q}) \in \mathbb{R}^{d_{\text{attn}}}, \quad (13)$$

where we omit the dependency on  $i$  for simplicity. Note that computation cost associated with Eq. (13) is  $2d_{\text{attn}} \cdot n_{\text{ctx}}$  for each head. Finally, output from all heads are concatenated followed by a linear map to project to  $d_{\text{model}}$ .

Note that Eq. (13) has the same form as FFN in Eq. (8) except for the choice of nonlinearity. Hence, following a similar strategy in obtaining Spark FFN, here we present Spark Attention as

$$\text{Spark-Attention}(\mathbf{q}; \mathbf{K}, \mathbf{V}, k, r) \stackrel{\text{def}}{=} \mathbf{V} \cdot \left( \text{softmax} \left( \text{Statistical-Top}_k^{(-\infty)}(\mathbf{K}^\top \mathbf{P}\mathbf{q}) \right) \odot \text{softplus}(\mathbf{K}^\top (\mathbf{I} - \mathbf{P})\mathbf{q}) \right) \quad (14)$$

In above,  $\mathbf{P} \stackrel{\text{def}}{=} \mathbf{1}_r \oplus \mathbf{0}_{d_{\text{attn}}-r} \in \mathbb{R}^{d_{\text{attn}} \times d_{\text{attn}}}$  is a fixed matrix.  $\text{Statistical-Top}_k^{(-\infty)}$  is a slight variant of Eq. (2) where the entries below the threshold  $\theta(\mathbf{x}, k)$  are set to  $-\infty$  instead of 0, so that such entries become zero after passing through softmax. Specifically,

$$[\text{Statistical-Top}_k^{(-\infty)}(\mathbf{x})]_i \stackrel{\text{def}}{=} \begin{cases} \mathbf{x}_i - \theta(\mathbf{x}, k) & \text{if } \mathbf{x}_i > \theta(\mathbf{x}, k), \\ -\infty & \text{otherwise.} \end{cases} \quad (15)$$

Finally, a softplus nonlinearity defined as the entrywise softplus function, i.e.,  $\log(1 + \exp(x))$ , is applied to the term  $\mathbf{K}^\top (\mathbf{I} - \mathbf{P})\mathbf{q}$  as this is empirically observed to offer quality benefits.

**FLOPs per Token.** With a naive implementation the number of FLOPs in Eq. (14) is given by  $2d_{\text{attn}} \cdot n_{\text{ctx}}$ , which is the same as the FLOPs for Eq. (13). However, by noting that the output of the softmax is expected to be sparse with approximately  $k$  nonzero entries, the computation costs associated with  $\mathbf{K}^\top (\mathbf{I} - \mathbf{P})\mathbf{q}$  and in the multiplication with  $\mathbf{V}$  can be drastically reduced. In particular, if we take  $r = \frac{d_{\text{attn}}}{2}$  then the FLOPs per token becomes

$$0.5d_{\text{model}}n_{\text{ctx}} + 1.5d_{\text{model}} \min\{k_{\text{attn}}, n_{\text{ctx}}\}, \quad (16)$$

which is nearly a  $4\times$  reduction when  $k_{\text{attn}}$  is much smaller than  $n_{\text{ctx}}$ .

## 4 EXPERIMENTS

In this section, we present an experimental evaluation of Spark Transformer using the Gemma-2 2B model. Gemma-2 2B is a decoder-only Transformer with 2 billion parameters, pretrained on 2 trillion tokens of primarily English text data (see Gemma Team (2024) for details). To evaluate Spark Transformer, we train a *Spark Gemma-2 2B* model by substituting the standard FFN and Attention in Gemma-2 2B with their Spark Transformer counterparts (Spark FFN and Spark Attention, respectively). This Spark Gemma-2 2B model is trained using the same procedure and data as the original Gemma-2 2B model.

**Implementation details.** Gemma-2 uses a model dimension of  $d_{\text{model}} = 2304$ . **For FFN**, Gemma-2 uses the Gated FFN in Eq. (12) with  $d_{\text{ff}}^g = 9216$ . We replace it with Spark FFN in Eq. (9) with  $d_{\text{ff}} = 13824$  so that the parameter count keeps the same. In addition, we take  $k = 1106$ , which gives a sparsity level of 8%, and  $r = 1024$  (due to sharding constraints,  $r$  can only be a multiple of 256). **For Attention**, Gemma-2 alternates between a global attention that have a span of 8192 tokens, and a local attention with a 4096 window size, both with  $d_{\text{attn}} = 256$ . We replace both with Spark Attention in Eq. (13) where for the latter we use the same 4096 window size. For hyper-parameters, we use  $k = 256$ , i.e. each token attends to at most 256 tokens, and  $r = 128$ . Extra care need to be taken for handling position embedding. Gemma-2 uses Rotary Position Embedding (Su et al., 2024) which is applied to  $\mathbf{q}$  and the columns of  $\mathbf{K}$  in Eq. (13). For Spark Attention in Eq. (14), we apply this position encoding to  $\mathbf{P}\mathbf{q}$ ,  $(\mathbf{I} - \mathbf{P})\mathbf{q}$ , the columns of  $\mathbf{P}\mathbf{K}$ , and the columns of  $(\mathbf{I} - \mathbf{P})\mathbf{K}$ .

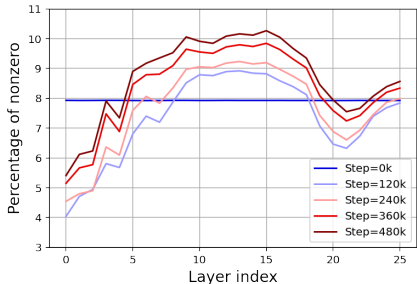
### 4.1 QUALITY

We evaluate Spark Gemma-2 2B on a suite of benchmarks that are used in the Gemma-2 paper (Gemma Team, 2024), and report the result in Table 2. We observe that Spark Gemma 2 matches the quality of Gemma 2 while having a drastically reduced FLOP count per token.

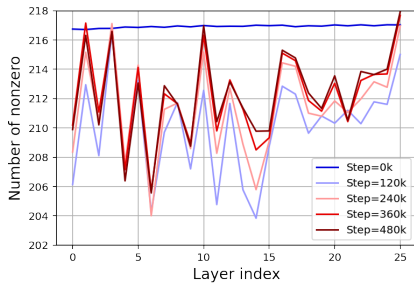
**Sparsity.** To verify the effectiveness of statistical top- $k$ , we report the level of sparsity measured in terms of percentage of nonzeros in FFN and the number of nonzeros in Attention. At the beginning of model training, we observe that statistical top- $k$  produces close to 8% nonzeros in FFN (see Figure 1a), which aligns well with our hyper-parameter choice of using  $k/d_{\text{ff}} = 8\%$  in Spark FFN. This

Table 2: Evaluation of Spark Transformer quality. We train a *Spark Gemma-2* model, replacing the standard FFN and Attention in Gemma-2 with Spark FFN and Spark Attention, respectively. We compare Spark Gemma-2 with ProSparse and LLaMA ReGLU, two recent models employing activation sparsity in their FFNs. Numbers in parentheses are taken from the respective original papers. FLOPs per token are computed assuming a context length of 8k.

|               | ProSparse<br>(Song et al., 2024a) | LLaMA ReGLU<br>(Zhang et al., 2024) | Gemma-2     | Spark Gemma-2<br>(Ours) |
|---------------|-----------------------------------|-------------------------------------|-------------|-------------------------|
| Model size    | 7B                                | 7B                                  | 2B          | 2B                      |
| FLOPs / token | -                                 | -                                   | 4.2B        | 1.4B                    |
| MMLU          | (45.5)                            | (44.8)                              | 52.1 (52.2) | 50.2                    |
| ARC-C         | -                                 | -                                   | 50.1 (55.7) | 51.1                    |
| GSM8K         | (12.1)                            | (10.6)                              | 21.2 (24.3) | 21.2                    |
| AGIEval       | (27.5)                            | -                                   | 31.8 (31.5) | 31.4                    |
| BBH           | (35.0)                            | -                                   | 41.3 (41.9) | 38.8                    |
| Winogrande    | -                                 | (69.4)                              | 68.7 (71.3) | 67.3                    |
| HellaSwag     | -                                 | (74.7)                              | 73.9 (72.9) | 73.2                    |
| MATH          | -                                 | -                                   | 16.4 (16.0) | 15.6                    |
| ARC-e         | -                                 | -                                   | 80.6 (80.6) | 81.3                    |
| PIQA          | -                                 | -                                   | 78.5 (78.4) | 78.6                    |
| SIQA          | -                                 | -                                   | 51.6 (51.9) | 51.3                    |
| Boolq         | -                                 | -                                   | 72.9 (72.7) | 73.3                    |
| TriviaQA      | -                                 | -                                   | 60.4 (60.4) | 59.4                    |
| NQ            | -                                 | -                                   | 17.1 (17.1) | 17.1                    |
| HumanEval     | -                                 | -                                   | 4.3 (20.1)  | 6.1                     |
| MBPP          | -                                 | -                                   | 30.4 (30.2) | 29.0                    |
| Avg.          | -                                 | -                                   | 47.0        | 46.6                    |



(a) Spark FFN sparsity



(b) Spark Attention sparsity

Figure 1: Sparsity in the intermediate activation of Spark FFN (i.e., output of GELU) and intermediate activation of Spark Attention (i.e., output of softmax) across all 26 layers at selected training steps. For FFN we report the percentage of nonzero entries out of  $d_{ff} = 13824$  entries. For Attention, we report the number of nonzero entries (i.e., attended tokens). Note that our hyper-parameter choice is to have 8% nonzeros in Spark FFN and at most 256 nonzeros in Spark Attention.

is expected as the model parameters, particularly  $K$  in Spark FFN, are randomly initialized, hence the entries of the activation maps are drawn from a Gaussian distribution which is in accordance with the assumption of statistical top- $k$ . The Gaussian assumption is no longer guaranteed after training, but we empirically observe it to hold approximately (see Section D.1) and statistical top- $k$  reliably produce a sparsity level close to 8% until the end of training at 480k steps. Sparsity in attention is reported in Figure 1b, which show that the number of attended tokens is below our hyper-parameter choice of 256 in Spark Attention throughout training. In particular, the numbers are much smaller because the results are from averaging over all tokens many of which have a context length of less than 256. Finally, we observe comparable levels of sparsity during evaluation (see Section D.2).

#### 4.2 EFFICIENCY

We evaluate the efficiency benefits of the Spark Gemma-2 2B over standard Gemma-2 2B using the `gemma.cpp` (Google Gemma.cpp, 2024), a C++ inference engine optimized for the Gemma models

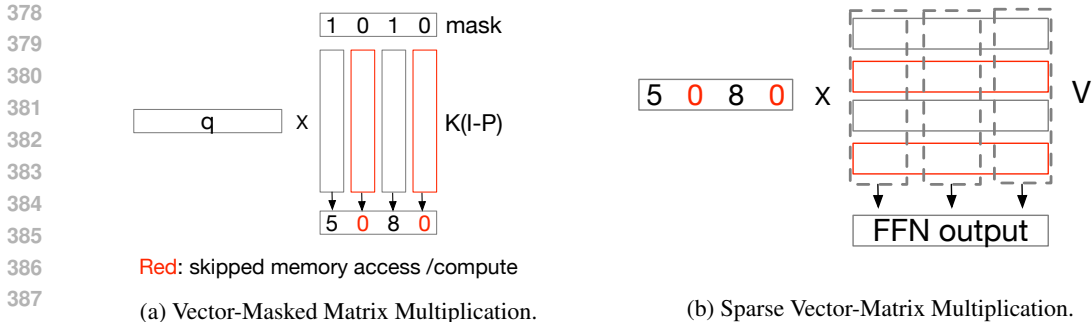


Figure 2: Illustration of the matrix multiplication implementation using sparse activation. (a) Vector-Masked Matrix Multiplication takes a dense vector  $q$ , a dense matrix  $K^T(I - P)$ , and a mask from statistical top- $k$  on  $K^T Pq$  to compute  $u := (K^T(I - P)q) \odot \text{mask}$ . It skips memory loading and compute associated with the masked columns. (b) Sparse Vector-Matrix Multiplication takes a sparse activation vector  $u$  to compute weighted sum of rows in the dense matrix  $V$ . It skips loading and computation of rows corresponding to 0's in  $u$ . To optimize performance, we implement Sparse Vector-Matrix Multiplication using tiling, which helps minimize cross-CPU core synchronization.

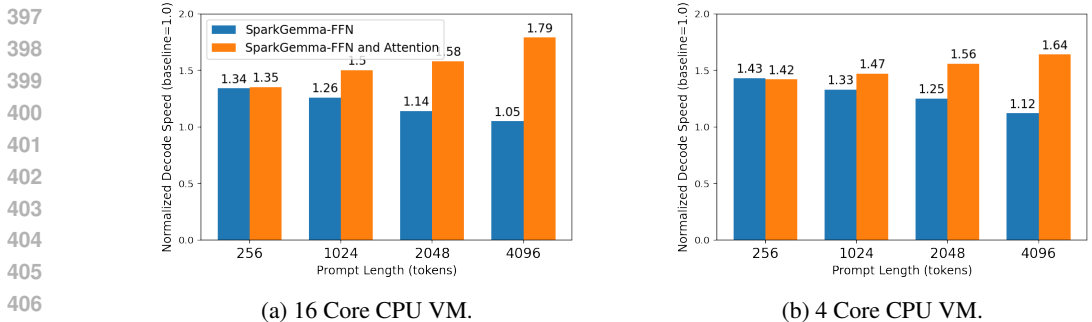


Figure 3: Spark Gemma-2 2B decoding speedup on CPU relative to the original Gemma-2 2B for varying prompt lengths. Speed is measured as the decoding time per token average over 128 tokens after the prompt. We provide a breakdown of speedup from FFN and Attention by reporting results of SparkGemma-FFN, which contains sparse optimization for FFN only, and SparkGemma-FFN and Attention, which contains sparse optimization for both. We use decode batch size of one.

on CPUs. Our implementation uses sparse matrix multiplication operators, which exploit sparsity in both FFN and Attention, as well as modern CPU vector SIMD operations (SIMD Wikipedia, 2024), see Figure 2 for an illustration and Section C in Appendix for details. We show that Spark Gemma-2 significantly improves the efficiency of transformer models, even in highly FLOP-constrained environments such as CPUs.

Specifically, Figure 3 reports the decoding speed under varying prompt lengths on a 4-Core or a 16-Core CPU. We see that Spark Gemma-2 outperforms the original Gemma-2 model, achieving a speedup that ranges from 1.35x to 1.79x on 16-Core CPU depending on the prompt length. For short prompts (e.g., 256 tokens), the sparse FFN optimization provides most of the speedup, whereas the sparse attention optimization provides the most speedup for longer prompts (e.g. 4096 tokens).

Table 3 further highlights the efficiency of Spark Gemma-2 in both prefill and decode phases. During the prefill, the prompt is usually chunked into batches since the process is bounded by memory bandwidth. This may reduce the benefit of activation sparsity as different tokens in a chunk may activate different subsets of parameters (in FFN) and attend to different subsets of tokens (in Attention). However, Table 3 shows that Spark Gemma-2 maintains strong performance with a chunk size of 64 tokens, following the default setup in `gemma.cpp`. A more detailed performance analysis of batching/chunking is provided in the Appendix. In addition, Spark Gemma-2 significantly outperforms the original Gemma-2 during decoding (with batch size=1). Notably, it achieves a decode speed of 86ms per token, which surpasses the average human reading speed (238 words per minute) (Brysbaert, 2019), with a very accessible 4-Core CPU Cloud VM.



Table 3: Prefill and decode speed of Spark Gemma-2 on 4-Core and 16-Core CPUs for prompts of 4096 tokens. During prefill phase, the prompt is chunked into batches of 64 tokens, following a default setup of `gemma.cpp`. Speedups relative to Gemma-2 are shown in parentheses.

|                | Prefill (ms / token) | Decode (ms / token) |
|----------------|----------------------|---------------------|
| 4 Core CPU VM  | 15 (1.86x)           | 86 (1.64x)          |
| 16 Core CPU VM | 7 (1.7x)             | 33 (1.79x)          |

### 4.3 ABLATION

**Statistical top- $k$ .** Adding a top- $k$  operator is expected to lead to a training slowdown due to the extra computation cost of computing top- $k$ . In Figure 4 we report the slowdown during training due to adding statistical top- $k$ . It can be observed that the slowdown is with a very small amount, demonstrating the efficiency of statistical top- $k$ . In particular, we compare the slowdown with that of the top- $k$  operator provided in JAX, namely `jax.lax.approx_max_k` (Chern et al., 2022). This operator is optimized to achieve TPU peak performance and has a controllable recall target, which we vary on the x-axis. Statistical top- $k$  is significantly faster than the JAX top- $k$  even when the latter operates on a small recall of 50%. Finally, we do not provide the quality of models trained with JAX top- $k$  since such models take a very long time to train.

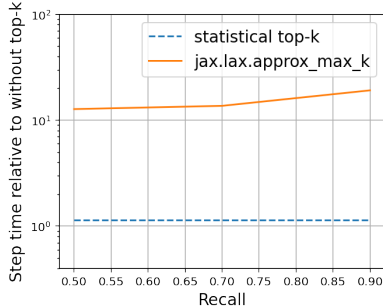


Figure 4: Comparison of training slowdown from using our statistical top- $k$  vs the standard top- $k$  (i.e., `jax.lax.approx_max_k` (Chern et al., 2022)) relative to not using any top- $k$ .

**Effect of  $r$  and  $k$  in Spark FFN.** Spark FFN comes with two hyper-parameters, namely  $r$  which controls the rank hence FLOP count of the low-cost predictor, and  $k$  which controls sparsity of activation hence the FLOP count. In Figure 5 we provide an ablation study on the effect of these two hyper-parameters, by reporting the training loss curves in the first 25,000 training steps (which is around 5% of full training). From Figure 5a, the best choice of  $r$  is 1024 which is nearly half of  $d_{\text{model}} = 2304$  (due to model sharding constraint,  $r$  cannot be taken to be exactly a half of  $d_{\text{model}}$ ). From Figure 5b, we see that the model quality is insensitive to choices of  $k$  that gives [5%, 10%] sparsity, but there is quality loss if we go sparser, e.g. 3% nonzeros.

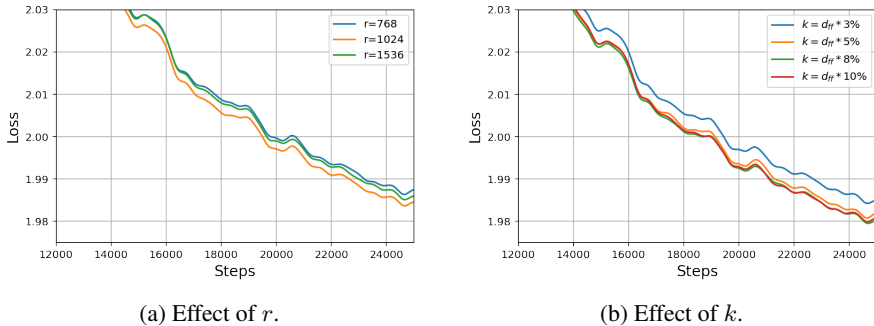


Figure 5: Effect of hyper-parameters  $r$  and  $k$  in Spark FFN on training loss. A Gaussian filter of  $\sigma = 200$  is applied to smooth the loss curves.

## 5 DISCUSSION AND RELATED WORK

Returning to the question posed at the outset: How many FLOPs is a token worth? This paper offers an answer through the Spark Transformer architecture, demonstrating at least a  $3\times$  overall FLOPs reduction without sacrificing model quality. This reduction is realized by selectively activating only part of the model parameters and limiting the attended context for each input. This principle of sparse activation finds a compelling parallel in neuroscience, where studies reveal sparse activity patterns in the brain as a key factor in its remarkable efficiency (Attwell & Laughlin, 2001; Barth & Poulet, 2012; Lee et al., 2024a). While hardware limitations currently hinder the full exploitation of

486 sparse activation in Transformers, particularly on GPUs and TPUs designed for dense computations,  
487 our work with Spark Transformer on CPUs highlights its potential. We believe this opens avenues  
488 for research into alternative hardware and platforms better suited for sparse computations, hence  
489 circumvents the hardware lottery (Hooker, 2021) and potentially lead to greater efficiency gains in  
490 the future.

491 In the following, we review a few lines of work closely related to ours.

492  
493 **Mixture of Experts (MoEs)** may be considered as a particular case of sparsely activated mod-  
494 els which group the neurons in FFN and activate all neurons in selected groups (Shazeer et al.,  
495 2017; Lepikhin et al., 2020). Neuron grouping has the benefit of being better suited for training  
496 accelerators compared to unstructured activation sparsity. However, training of MoEs incurs extra  
497 complexities in algorithmic design and requires special hardware support (Fedus et al., 2022). More-  
498 over, the structured nature of sparsity limits the model’s flexibility and expressiveness, and recent  
499 work advocates the use of a larger number of smaller experts (Dai et al., 2024; He, 2024). On the  
500 other hand, the discovery of the naturally emerging unstructured activation sparsity has motivated  
501 the new perspective of naturally emerging experts (Zhang et al., 2022; Dong et al., 2023; Csordás  
502 et al., 2023; Qiu et al., 2023; Szatkowski et al., 2024; Zheng et al., 2024).

503 **Sparse activation** is common approach to improve the efficiency of large models and many tech-  
504 niques for a low-cost activation prediction have been developed over the years, such as low-rank  
505 factorization (Davis & Arel, 2013), quantization (Cao et al., 2019), product keys (Lample et al.,  
506 2019), hashing (Chen et al., 2020), etc. With the popularity of modern large Transformer models,  
507 these techniques become natural choices (Jaszczur et al., 2021; Zeng et al., 2023; Liu et al., 2023;  
508 Song et al., 2023) for reducing their high computation costs. In particular, a lot of the excitement  
509 comes from the discovery that the activations in FFNs are naturally sparse (Zhang et al., 2022; Li  
510 et al., 2022) and hence efficiency with activation sparsity are obtained without a quality toll.

511 Our work falls into the category of the latest work in this direction that aim to bring the benefits to the  
512 latest generation large language models that do not have natural sparsity. Early attempts (Mirzadeh  
513 et al., 2023; Peng et al., 2023; Zhang et al., 2024) seek to bring back sparsity by switching back to  
514 ReLU variants, but it usually incurs a quality loss. The quality gap may be largely bridged by more  
515 careful tuning, but the activation becomes less sparse (e.g. 25% nonzeros in LLAMA 7B (Song  
516 et al., 2024a)). Top-k becomes a more popular choice for obtaining sparsity recently (Song et al.,  
517 2024b) and is able to maintain neutral quality while offering strong sparsity, but only in selected  
518 layers (Yerram et al., 2024). Moreover, such methods require finetuning to bring sparsity and also  
519 obtain a predictor. Without doing finetuning, Lee et al. (2024b); Liu et al. (2024a) obtained at most  
520 50% nonzeros under neutral quality. In contrast to these works, our work not only obtains 8%  
521 nonzeros in activation of all FFN layers, but also a predictor, all with a single-stage training. **We  
provide a summary of comparison to these methods in Table 4 in the Appendix.**

522 Finally, the usefulness of activation sparsity goes beyond efficiency. For example, theoretical studies  
523 show its benefits for model generalizability and learnability (Muthukumar & Sulam, 2023; Awasthi  
524 et al., 2024). Moreover, activated neurons may be associated with semantic concepts, which of-  
525 fers understanding of the working mechanism and enables manipulating the output of Transformer  
526 models (Cuadros et al., 2022; Luo et al., 2024).

527 **Sparse attention** broadly refers to the approach of attending to a selected subset of tokens in the con-  
528 text as a means of reducing computation cost (Deng et al., 2024; Jiang et al., 2024). Work on sparse  
529 attention include those that use handcrafted attention patterns (Child et al., 2019; Beltagy et al.,  
530 2020; Ainslie et al., 2023; Ding et al., 2023), which feature simplicity, and learned attention patterns  
531 (Kitaev et al., 2020; Roy et al., 2021) which feature better modeling capacity. However, learning  
532 attention patterns often involve learning, e.g., a hash table or k-means centers, which significantly  
533 complicates modeling. Closely related to our Spark Attention is the top- $k$  attention (Gupta et al.,  
534 2021), which obtains data-adaptive attention simply from top- $k$  thresholding. Our work improves  
535 upon top- $k$  attention by introducing a low cost predictor which enables an increased computational  
536 benefits from sparsity. **Finally, KV pruning approaches drop selected tokens permanently as decod-  
ing proceeds (Zhang et al., 2023; Liu et al., 2024b), and cannot achieve as high compression ratio as  
537 Top- $k$  based approaches.**

## REFERENCES

- 540  
541  
542 Bo Adler, Niket Agarwal, Ashwath Aithal, Dong H Anh, Pallab Bhattacharya, Annika Brundyn,  
543 Jared Casper, Bryan Catanzaro, Sharon Clay, Jonathan Cohen, et al. Nemotron-4 340b technical  
544 report. *arXiv preprint arXiv:2406.11704*, 2024.
- 545 Joshua Ainslie, Tao Lei, Michiel de Jong, Santiago Ontanon, Siddhartha Brahma, Yury Zemlyan-  
546 ski, David Uthus, Mandy Guo, James Lee-Thorp, Yi Tay, et al. Colt5: Faster long-range trans-  
547 formers with conditional computation. In *Conference on Empirical Methods in Natural Language*  
548 *Processing*, 2023.
- 549 Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Co-  
550 jocar, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic,  
551 et al. The falcon series of open language models. *arXiv preprint arXiv:2311.16867*, 2023.
- 553 Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos,  
554 Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report.  
555 *arXiv preprint arXiv:2305.10403*, 2023.
- 556 David Attwell and Simon B Laughlin. An energy budget for signaling in the grey matter of the  
557 brain. *Journal of Cerebral Blood Flow & Metabolism*, 21(10):1133–1145, 2001.
- 559 Pranjal Awasthi, Nishanth Dikkala, Pritish Kamath, and Raghu Meka. Learning neural networks  
560 with sparse activations. In *The Thirty Seventh Annual Conference on Learning Theory*, pp. 406–  
561 425. PMLR, 2024.
- 562 Alison L Barth and James FA Poulet. Experimental evidence for sparse firing in the neocortex.  
563 *Trends in neurosciences*, 35(6):345–355, 2012.
- 565 Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse  
566 problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.
- 567 Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer.  
568 *arXiv preprint arXiv:2004.05150*, 2020.
- 570 Thomas Blumensath and Mike E Davies. Iterative thresholding for sparse approximations. *Journal*  
571 *of Fourier analysis and Applications*, 14:629–654, 2008.
- 572 Marc Brysbaert. How many words do we read per minute? a review and meta-analysis of reading  
573 rate. *Journal of Memory and Language*, 109:104047, 2019.
- 575 Shijie Cao, Lingxiao Ma, Wencong Xiao, Chen Zhang, Yunxin Liu, Lintao Zhang, Lanshun Nie, and  
576 Zhi Yang. Seernet: Predicting convolutional neural network feature-map sparsity through low-  
577 bit quantization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
578 *Recognition*, pp. 11216–11225, 2019.
- 579 Beidi Chen, Tharun Medini, James Farwell, Charlie Tai, Anshumali Shrivastava, et al. Slide: In  
580 defense of smart algorithms over hardware acceleration for large-scale deep learning systems.  
581 *Proceedings of Machine Learning and Systems*, 2:291–306, 2020.
- 583 Felix Chern, Blake Hechtman, Andy Davis, Ruiqi Guo, David Majnemer, and Sanjiv Kumar. Tpu-  
584 knn: K nearest neighbor search at peak flop/s. *Advances in Neural Information Processing Sys-*  
585 *tems*, 35:15489–15501, 2022.
- 586 Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse  
587 transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- 589 Róbert Csordás, Kazuki Irie, and Jürgen Schmidhuber. Approximating two-layer feedforward net-  
590 works for efficient transformers. In *Findings of the Association for Computational Linguistics:*  
591 *EMNLP 2023*, pp. 674–692, 2023.
- 592 Xavier Suau Cuadros, Luca Zappella, and Nicholas Apostoloff. Self-conditioning pre-trained lan-  
593 guage models. In *International Conference on Machine Learning*, pp. 4455–4473. PMLR, 2022.

- 594 Damai Dai, Chengqi Deng, Chenggang Zhao, RX Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding  
595 Zeng, Xingkai Yu, Y Wu, et al. Deepseekmoe: Towards ultimate expert specialization in mixture-  
596 of-experts language models. *arXiv preprint arXiv:2401.06066*, 2024.
- 597  
598 Yann N Dauphin, Angela Fan, Michael Auli, and David Grangier. Language modeling with gated  
599 convolutional networks. In *International conference on machine learning*, pp. 933–941. PMLR,  
600 2017.
- 601 Andrew Davis and Itamar Arel. Low-rank approximations for conditional feedforward computation  
602 in deep neural networks. *arXiv preprint arXiv:1312.4461*, 2013.
- 603  
604 Yichuan Deng, Zhao Song, and Chiwun Yang. Attention is naturally sparse with gaussian distributed  
605 input. *arXiv preprint arXiv:2404.02690*, 2024.
- 606  
607 Jiayu Ding, Shuming Ma, Li Dong, Xingxing Zhang, Shaohan Huang, Wenhui Wang, Nanning  
608 Zheng, and Furu Wei. Longnet: Scaling transformers to 1,000,000,000 tokens. *arXiv preprint*  
609 *arXiv:2307.02486*, 2023.
- 610 Harry Dong, Beidi Chen, and Yuejie Chi. Towards structured sparsity in transformers for efficient  
611 inference. In *Workshop on Efficient Systems for Foundation Models@ ICML2023*, 2023.
- 612  
613 Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale.  
614 *arXiv preprint arXiv:2010.11929*, 2020.
- 615  
616 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
617 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
618 *arXiv preprint arXiv:2407.21783*, 2024.
- 619 Aryeh Dvoretzky, Jack Kiefer, and Jacob Wolfowitz. Asymptotic minimax character of the sample  
620 distribution function and of the classical multinomial estimator. *The Annals of Mathematical*  
621 *Statistics*, pp. 642–669, 1956.
- 622  
623 William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter  
624 models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39,  
625 2022.
- 626  
627 Gemma Team. Gemma 2: Improving open language models at a practical size. *arXiv preprint*  
628 *arXiv:2408.00118*, 2024.
- 629  
630 Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are  
631 key-value memories. In *Proceedings of the 2021 Conference on Empirical Methods in Natural*  
*Language Processing*, pp. 5484–5495, 2021.
- 632  
633 Google Gemma.cpp. Google/gemma.cpp: lightweight, standalone c++ inference engine for google’s  
634 gemma models. <https://github.com/google/gemma.cpp>, 2024.
- 635  
636 Ankit Gupta, Guy Dar, Shaya Goodman, David Ciprut, and Jonathan Berant. Memory-efficient  
637 transformers via top- $k$  attention. *arXiv preprint arXiv:2106.06899*, 2021.
- 638  
639 Xu Owen He. Mixture of a million experts. *arXiv preprint arXiv:2407.04153*, 2024.
- 640  
641 Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint*  
642 *arXiv:1606.08415*, 2016.
- 643  
644 Sara Hooker. The hardware lottery. *Communications of the ACM*, 64(12):58–65, 2021.
- 645  
646 Peter J Huber. Robust estimation of a location parameter. In *Breakthroughs in statistics: Methodol-*  
647 *ogy and distribution*, pp. 492–518. Springer, 1992.
- 648  
649 Sebastian Jaszczur, Aakanksha Chowdhery, Afroz Mohiuddin, Lukasz Kaiser, Wojciech Gajewski,  
650 Henryk Michalewski, and Jonni Kanerva. Sparse is enough in scaling transformers. *Advances in*  
*Neural Information Processing Systems*, 34:9895–9907, 2021.

- 648 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
649 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.  
650 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.  
651
- 652 Huiqiang Jiang, YUCHENG LI, Chengruidong Zhang, Qianhui Wu, Xufang Luo, Surin Ahn, Zhen-  
653 hua Han, Amir H Abdi, Dongsheng Li, Chin-Yew Lin, et al. Minference 1.0: Accelerating pre-  
654 filling for long-context llms via dynamic sparse attention. In *The Thirty-eighth Annual Conference*  
655 *on Neural Information Processing Systems*, 2024.
- 656 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,  
657 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language  
658 models. *arXiv preprint arXiv:2001.08361*, 2020.  
659
- 660 Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In  
661 *International Conference on Learning Representations*, 2020.  
662
- 663 Guillaume Lample, Alexandre Sablayrolles, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé  
664 Jégou. Large memory layers with product keys. *Advances in Neural Information Processing*  
665 *Systems*, 32, 2019.
- 666 Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selec-  
667 tion. *Annals of statistics*, pp. 1302–1338, 2000.  
668
- 669 J Quinn Lee, Matt Nielsen, Rebecca McHugh, Erik Morgan, Nancy S Hong, Robert J Sutherland,  
670 and Robert J McDonald. Sparsity of population activity in the hippocampus is task-invariant  
671 across the trisynaptic circuit and dorsoventral axis. *Behavioural Brain Research*, 423:113790,  
672 2024a.
- 673 Je-Yong Lee, Donghyun Lee, Genghan Zhang, Mo Tiwari, and Azalia Mirhoseini. Cats:  
674 Contextually-aware thresholding for sparsity in large language models. *arXiv preprint*  
675 *arXiv:2404.08763*, 2024b.  
676
- 677 Wonbeom Lee, Jungi Lee, Junghwan Seo, and Jaewoong Sim. {InfiniGen}: Efficient generative  
678 inference of large language models with dynamic {KV} cache management. In *18th USENIX*  
679 *Symposium on Operating Systems Design and Implementation (OSDI 24)*, pp. 155–172, 2024c.
- 680 Tao Lei, Junwen Bai, Siddhartha Brahma, Joshua Ainslie, Kenton Lee, Yanqi Zhou, Nan Du, Vincent  
681 Zhao, Yuexin Wu, Bo Li, et al. Conditional adapters: Parameter-efficient transfer learning with  
682 fast inference. *Advances in Neural Information Processing Systems*, 36:8152–8172, 2023.  
683
- 684 Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang,  
685 Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional  
686 computation and automatic sharding. In *International Conference on Learning Representations*,  
687 2020.
- 688 Zonglin Li, Chong You, Srinadh Bhojanapalli, Daliang Li, Ankit Singh Rawat, Sashank J Reddi,  
689 Ke Ye, Felix Chern, Felix Yu, Ruiqi Guo, et al. The lazy neuron phenomenon: On emergence  
690 of activation sparsity in transformers. In *The Eleventh International Conference on Learning*  
691 *Representations*, 2022.  
692
- 693 Yujun Lin, Song Han, Huizi Mao, Yu Wang, and Bill Dally. Deep gradient compression: Reducing  
694 the communication bandwidth for distributed training. In *International Conference on Learning*  
695 *Representations*, 2018.
- 696 James Liu, Pragaash Ponnusamy, Tianle Cai, Han Guo, Yoon Kim, and Ben Athiwaratkun. Training-  
697 free activation sparsity in large language models. *arXiv preprint arXiv:2408.14690*, 2024a.  
698
- 699 Zichang Liu, Jue Wang, Tri Dao, Tianyi Zhou, Binhang Yuan, Zhao Song, Anshumali Shrivastava,  
700 Ce Zhang, Yuandong Tian, Christopher Re, et al. Deja vu: Contextual sparsity for efficient llms  
701 at inference time. In *International Conference on Machine Learning*, pp. 22137–22176. PMLR,  
2023.

- 702 Zichang Liu, Aditya Desai, Fangshuo Liao, Weitao Wang, Victor Xie, Zhaozhuo Xu, Anastasios  
703 Kyrillidis, and Anshumali Shrivastava. Scissorhands: Exploiting the persistence of importance  
704 hypothesis for llm kv cache compression at test time. *Advances in Neural Information Processing*  
705 *Systems*, 36, 2024b.
- 706  
707 Chao Lou, Zixia Jia, Zilong Zheng, and Kewei Tu. Sparser is faster and less is more: Efficient sparse  
708 attention for long-range transformers. *arXiv preprint arXiv:2406.16747*, 2024.
- 709  
710 Jinqi Luo, Tianjiao Ding, Kwan Ho Ryan Chan, Darshan Thaker, Aditya Chattopadhyay, Chris  
711 Callison-Burch, and René Vidal. Pace: Parsimonious concept engineering for large language  
712 models. *arXiv preprint arXiv:2406.04331*, 2024.
- 713  
714 Ahmed M Abdelmoniem, Ahmed Elzanaty, Mohamed-Slim Alouini, and Marco Canini. An efficient  
715 statistical-based gradient compression technique for distributed training systems. *Proceedings of*  
716 *Machine Learning and Systems*, 3:297–322, 2021.
- 717  
718 Pascal Massart. The tight constant in the dvoretzky-kiefer-wolfowitz inequality. *The annals of*  
719 *Probability*, pp. 1269–1283, 1990.
- 720  
721 Iman Mirzadeh, Keivan Alizadeh, Sachin Mehta, Carlo C Del Mundo, Oncel Tuzel, Golnoosh  
722 Samei, Mohammad Rastegari, and Mehrdad Farajtabar. Relu strikes back: Exploiting activation  
723 sparsity in large language models. *arXiv preprint arXiv:2310.04564*, 2023.
- 724  
725 Ramchandran Muthukumar and Jeremias Sulam. Sparsity-aware generalization theory for deep  
726 neural networks. In *Annual Conference on Learning Theory*, pp. 5311–5342. PMLR, 2023.
- 727  
728 Neal Parikh, Stephen Boyd, et al. Proximal algorithms. *Foundations and trends® in Optimization*,  
729 1(3):127–239, 2014.
- 730  
731 David Patterson, Joseph Gonzalez, Quoc Le, Chen Liang, Lluís-Miquel Munguia, Daniel Rothchild,  
732 David So, Maud Texier, and Jeff Dean. Carbon emissions and large neural network training. *arXiv*  
733 *preprint arXiv:2104.10350*, 2021.
- 734  
735 Ze Peng, Lei Qi, Yinghuan Shi, and Yang Gao. Theoretical explanation of activation sparsity through  
736 flat minima and adversarial robustness. *arXiv preprint arXiv:2309.03004*, 2023.
- 737  
738 Zihan Qiu, Zeyu Huang, and Jie Fu. Emergent mixture-of-experts: Can dense pre-trained transform-  
739 ers benefit from emergent modular structures? *arXiv preprint arXiv:2310.10908*, 2023.
- 740  
741 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
742 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text  
743 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 744  
745 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-  
746 baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gem-  
747 ini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint*  
748 *arXiv:2403.05530*, 2024.
- 749  
750 Luka Ribar, Ivan Chelombiev, Luke Hudliss-Galley, Charlie Blake, Carlo Luschi, and Douglas Orr.  
751 Sparq attention: Bandwidth-efficient llm inference. *arXiv preprint arXiv:2312.04985*, 2023.
- 752  
753 Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. Efficient content-based sparse  
754 attention with routing transformers. *Transactions of the Association for Computational Linguis-*  
755 *tics*, 9:53–68, 2021.
- Or Sharir, Barak Peleg, and Yoav Shoham. The cost of training nlp models: A concise overview.  
*arXiv preprint arXiv:2004.08900*, 2020.
- Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton,  
and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer.  
*arXiv preprint arXiv:1701.06538*, 2017.

- 756 Shaohuai Shi, Xiaowen Chu, Ka Chun Cheung, and Simon See. Understanding top-k sparsification  
757 in distributed deep learning. *arXiv preprint arXiv:1911.08772*, 2019.  
758
- 759 SIMD Wikipedia. Single instruction, multiple data - wikipedia — en.wikipedia.org. [https://](https://en.wikipedia.org/wiki/Single_instruction,_multiple_data)  
760 [en.wikipedia.org/wiki/Single\\_instruction,\\_multiple\\_data](https://en.wikipedia.org/wiki/Single_instruction,_multiple_data), 2024.  
761
- 762 Chenyang Song, Xu Han, Zhengyan Zhang, Shengding Hu, Xiyu Shi, Kuai Li, Chen Chen, Zhiyuan  
763 Liu, Guangli Li, Tao Yang, et al. Prospare: Introducing and enhancing intrinsic activation spar-  
764 sity within large language models. *arXiv preprint arXiv:2402.13516*, 2024a.
- 765 Yixin Song, Zeyu Mi, Haotong Xie, and Haibo Chen. Powerinfer: Fast large language model serving  
766 with a consumer-grade gpu. *arXiv preprint arXiv:2312.12456*, 2023.  
767
- 768 Yixin Song, Haotong Xie, Zhengyan Zhang, Bo Wen, Li Ma, Zeyu Mi, and Haibo Chen. Turbo  
769 sparse: Achieving llm sota performance with minimal activated parameters. *arXiv preprint*  
770 *arXiv:2406.05955*, 2024b.
- 771 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: En-  
772 hanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.  
773
- 774 Filip Szatkowski, Bartosz Wójcik, Mikołaj Piórczyński, and Simone Scardapane. Exploiting acti-  
775 vation sparsity with dense to dynamic-k mixture-of-experts conversion. In *Workshop on Efficient*  
776 *Systems for Foundation Models II@ ICML2024*, 2024.  
777
- 778 Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau,  
779 Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der  
780 Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson,  
781 Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W. Moore,  
782 Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero,  
783 Charles R. Harris, Anne M. Archibald, António H. Ribeiro, Fabian Pedregosa, Paul van Mul-  
784 bregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing  
785 in Python. *Nature Methods*, 17:261–272, 2020. doi: 10.1038/s41592-019-0686-2.
- 786 Shuo Yang, Ying Sheng, Joseph E Gonzalez, Ion Stoica, and Lianmin Zheng. Post-training sparse  
787 attention with double sparsity. *arXiv preprint arXiv:2408.07092*, 2024.  
788
- 789 Varun Yerram, Chong You, Srinadh Bhojanapalli, Sanjiv Kumar, Prateek Jain, Praneeth Netrapalli,  
790 et al. Hire: High recall approximate top- $k$  estimation for efficient llm inference. *arXiv preprint*  
791 *arXiv:2402.09360*, 2024.  
792
- 793 Zhanpeng Zeng, Michael Davies, Pranav Pulijala, Karthikeyan Sankaralingam, and Vikas Singh.  
794 Lookuppfn: making transformers compute-lite for cpu inference. In *International Conference on*  
795 *Machine Learning*, pp. 40707–40718. PMLR, 2023.
- 796 Zhengyan Zhang, Yankai Lin, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie Zhou. Moefication:  
797 Transformer feed-forward layers are mixtures of experts. In *Findings of the Association for Com-*  
798 *putational Linguistics: ACL 2022*, pp. 877–890, 2022.  
799
- 800 Zhengyan Zhang, Yixin Song, Guanghui Yu, Xu Han, Yankai Lin, Chaojun Xiao, Chenyang Song,  
801 Zhiyuan Liu, Zeyu Mi, and Maosong Sun. Relu<sup>2</sup> wins: Discovering efficient activation functions  
802 for sparse llms. *arXiv preprint arXiv:2402.03804*, 2024.  
803
- 804 Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song,  
805 Yuandong Tian, Christopher Ré, Clark Barrett, et al. H2o: Heavy-hitter oracle for efficient gen-  
806 erative inference of large language models. *Advances in Neural Information Processing Systems*,  
807 36:34661–34710, 2023.
- 808 Haizhong Zheng, Xiaoyan Bai, Beidi Chen, Fan Lai, and Atul Prakash. Learn to be efficient: Build  
809 structured sparsity in large language models. *arXiv preprint arXiv:2402.06126*, 2024.

## A PROOF TO THEOREM 1

*Proof.* In this proof we write  $\bar{x} \stackrel{\text{def}}{=} \text{mean}(x)$  and  $s \stackrel{\text{def}}{=} \text{std}(x)$  for brevity.

We first establish the concentration bounds that the empirical mean and standard deviation, i.e.,  $\bar{x}$  and  $s$  are close to the true mean and true standard deviation, i.e.,  $\mu$  and  $\sigma$  of the underlying Gaussian, respectively. Recall from the definition of the chi-squared distribution that  $(d-1)\frac{s^2}{\sigma^2} \sim \chi^2(d-1)$ . Using the Laurent-Massart bound on the tail probability of the chi-squared distribution (Laurent & Massart, 2000, Corollary of Lemma 1), we have

$$\Pr\left(\left|(d-1)\frac{s^2}{\sigma^2} - (d-1)\right| \geq 2\sqrt{(d-1)t} + 2t\right) \leq 2e^{-t}$$

for every  $t > 0$ . We set  $t = \log \frac{6}{\delta}$ . Then, with a probability of at least  $1 - \delta/3$ , we have

$$(d-1)\left|\frac{s^2}{\sigma^2} - 1\right| < 2\sqrt{(d-1)\log \frac{6}{\delta}} + 2\log \frac{6}{\delta},$$

which implies

$$\left|\frac{s^2}{\sigma^2} - 1\right| < 2\sqrt{\frac{\log \frac{6}{\delta}}{d-1}} + 2\frac{\log \frac{6}{\delta}}{d-1} \leq 4\sqrt{\frac{\log \frac{6}{\delta}}{d}} + 4\frac{\log \frac{6}{\delta}}{d} \leq 8\sqrt{\frac{\log \frac{6}{\delta}}{d}},$$

where the last inequality uses the assumption that  $d \geq \max\{2, \log \frac{6}{\delta}\}$ . By rearranging the terms, we get

$$\sigma\left(1 - 8\sqrt{\frac{\log \frac{6}{\delta}}{d}}\right) \leq \sigma\sqrt{\max\left\{1 - 8\sqrt{\frac{\log \frac{6}{\delta}}{d}}, 0\right\}} \leq s \leq \sigma\sqrt{1 + 8\sqrt{\frac{\log \frac{6}{\delta}}{d}}} \leq \sigma\left(1 + 8\sqrt{\frac{\log \frac{6}{\delta}}{d}}\right),$$

which simplifies to

$$|s - \sigma| \leq 8\sigma\sqrt{\frac{\log \frac{6}{\delta}}{d}}. \quad (17)$$

Eq. (17) provides a concentration bound for  $s$ . We now proceed to deriving a bound for  $\mu$ . Towards that, notice that  $\frac{\bar{x} - \mu}{\sigma/\sqrt{d}} \sim \mathcal{N}(0, 1)$ . By using the Mill's inequality that upper bounds the tail probability of a standard normal distribution (i.e., if  $Z \sim \mathcal{N}(0, 1)$  and  $t > 0$ , then  $\Pr(|Z| > t) \leq \frac{e^{-t^2/2}}{t}$ ), we have

$$\Pr\left(\left|\frac{\bar{x} - \mu}{\sigma/\sqrt{d}}\right| > \sqrt{2\log \frac{3}{\delta}}\right) \leq \frac{\delta/3}{\sqrt{2\log \frac{3}{\delta}}} \leq \delta/3.$$

Therefore, with probability at least  $1 - \delta/3$ , we have

$$\left|\frac{\bar{x} - \mu}{\sigma/\sqrt{d}}\right| \leq \sqrt{2\log \frac{3}{\delta}},$$

which yields

$$|\bar{x} - \mu| \leq \sigma\sqrt{\frac{2\log \frac{3}{\delta}}{d}}. \quad (18)$$

Combining Eq. (17) and Eq. (18), with probability at least  $1 - 2\delta/3$ , we have

$$\left|\theta(\mathbf{x}, k) - \left(\mu + \sigma Q\left(1 - \frac{k}{d}\right)\right)\right| \quad (19)$$

$$\leq |\bar{x} - \mu| + |s - \sigma| \left|Q\left(1 - \frac{k}{d}\right)\right| \quad (20)$$

$$\leq \sigma\sqrt{\frac{2\log \frac{3}{\delta}}{d}} + 8\sigma\sqrt{\frac{\log \frac{6}{\delta}}{d}} \left|Q\left(1 - \frac{k}{d}\right)\right|. \quad (21)$$



We define the empirical cumulative distribution function (ECDF) of  $x_1, x_2, \dots, x_d$  as

$$\hat{F}_d(x) = \frac{1}{d} \sum_{i \in [d]} \mathbf{1}_{\{x_i \leq x\}}.$$

Then, the number of the entries of  $\mathbf{x}$  that are greater than  $\theta(\mathbf{x}, k)$  may be written as

$$\text{card}(\{i \in [d] \mid x_i > \theta(\mathbf{x}, k)\}) = \sum_{i \in [d]} \mathbf{1}_{\{x_i > \theta(\mathbf{x}, k)\}} = d(1 - \hat{F}_d(\theta(\mathbf{x}, k))).$$

Let  $F$  denote the cumulative distribution function (CDF) of  $\mathcal{N}(\mu, \sigma^2)$ . By the Dvoretzky-Kiefer-Wolfowitz inequality (Dvoretzky et al., 1956; Massart, 1990), we have

$$\Pr\left(\sup_{u \in \mathbb{R}} |\hat{F}_d(u) - F(u)| > t\right) \leq 2e^{-2dt^2}.$$

Taking  $t = \sqrt{\frac{1}{2d} \log \frac{6}{\delta}}$  and  $u = \theta(\mathbf{x}, k)$ , we obtain

$$\Pr\left(\left|\hat{F}_d(\theta(\mathbf{x}, k)) - F(\theta(\mathbf{x}, k))\right| > \sqrt{\frac{1}{2d} \log \frac{6}{\delta}}\right) \leq \frac{\delta}{3}. \quad (22)$$

Applying the union bound on Eq. (19) and Eq. (22), we obtain that the following holds with probability at least  $1 - \delta$ :

$$\left|\hat{F}_d(\theta(\mathbf{x}, k)) - \left(1 - \frac{k}{d}\right)\right| \quad (23)$$

$$= \left|\hat{F}_d(\theta(\mathbf{x}, k)) - F(\mu + \sigma Q(1 - \frac{k}{d}))\right| \quad (24)$$

$$\leq \left|\hat{F}_d(\theta(\mathbf{x}, k)) - F(\theta(\mathbf{x}, k))\right| + \left|F(\theta(\mathbf{x}, k)) - F(\mu + \sigma Q(1 - \frac{k}{d}))\right| \quad (25)$$

$$\leq \sqrt{\frac{1}{2d} \log \frac{6}{\delta}} + \frac{1}{\sqrt{2\pi}\sigma} \left|\theta(\mathbf{x}, k) - (\mu + \sigma Q(1 - \frac{k}{d}))\right| \quad (26)$$

$$\leq \sqrt{\frac{1}{2d} \log \frac{6}{\delta}} + \frac{1}{\sqrt{2\pi}} \left(\sqrt{\frac{2 \log \frac{3}{\delta}}{d}} + 8\sqrt{\frac{\log \frac{6}{\delta}}{d}} \left|Q(1 - \frac{k}{d})\right|\right) \quad (27)$$

$$\leq 4\sqrt{\frac{\log \frac{6}{\delta}}{d}} \left(1 + \left|Q(1 - \frac{k}{d})\right|\right). \quad (28)$$

In the above expression, the first equality stems directly from the definitions of  $F(\cdot)$  and  $Q(\cdot)$ , which gives

$$F(\theta(\mathbf{x}, k)) = F(\mu + \sigma Q(1 - \frac{k}{d})) = \Phi(Q(1 - \frac{k}{d})) = 1 - \frac{k}{d},$$

where  $\Phi$  denotes the CDF of the standard normal distribution.

To simplify Eq. (23), we consider two cases:

- If  $k \leq d/2$ , by Mill's inequality, we have

$$1 - \Phi\left(\sqrt{2 \log \frac{d}{k}}\right) \leq \frac{e^{-(\sqrt{2 \log \frac{d}{k}})^2/2}}{\sqrt{2 \log \frac{d}{k}}} = \frac{e^{-(\sqrt{2 \log \frac{d}{k}})^2/2}}{\sqrt{2 \log \frac{d}{k}}} = \frac{k/d}{\sqrt{2 \log \frac{d}{k}}} \leq \frac{k}{d},$$

where the last inequality is because  $2 \log \frac{d}{k} \geq 1$ . Therefore

$$1 - \frac{k}{d} \leq \Phi\left(\sqrt{2 \log \frac{d}{k}}\right),$$

which gives

$$Q(1 - \frac{k}{d}) \leq \sqrt{2 \log \frac{d}{k}} = \sqrt{-2 \log \frac{k}{d}}.$$

- If  $k > d/2$ , we have

$$\left|Q\left(1 - \frac{k}{d}\right)\right| = \left|Q\left(1 - \frac{d-k}{d}\right)\right| \leq \sqrt{-2 \log \frac{d-k}{d}}.$$

Combining the two cases, we get

$$\left|Q\left(1 - \frac{k}{d}\right)\right| \leq \sqrt{-2 \log \min \left\{\frac{k}{d}, 1 - \frac{k}{d}\right\}}.$$

Plugging this into Eq. (23), we obtain

$$\left|\hat{F}_d(\theta(\mathbf{x}, k)) - \left(1 - \frac{k}{d}\right)\right| \leq 4\sqrt{\frac{\log \frac{6}{\delta}}{d}} \left(1 + \sqrt{-2 \log \min \left\{\frac{k}{d}, 1 - \frac{k}{d}\right\}}\right).$$

Recall  $\text{card}(\{i \in [d] \mid x_i > \theta(\mathbf{x}, k)\}) = d\left(1 - \hat{F}_d(\theta(\mathbf{x}, k))\right)$ . We conclude that with probability at least  $1 - \delta$ , we have

$$\left|\text{card}(\{i \in [d] \mid x_i > \theta(\mathbf{x}, k)\}) - k\right| \leq 4\sqrt{d \log \frac{6}{\delta}} \left(1 + \sqrt{-2 \log \min \left\{\frac{k}{d}, 1 - \frac{k}{d}\right\}}\right).$$

□

## B PROOF TO THEOREM 2

*Proof.* The Huber statistical top- $k$  in Eq. (4) may be written as

$$\text{Huber}(\text{Statistical-Top}_k(\mathbf{x}); \delta) / \delta = \text{Huber}(\text{Soft-Threshold}(\mathbf{x}, \theta(\mathbf{x}, k))) / \delta, \quad (29)$$

where  $\theta(\mathbf{x}, k)$  is defined in Eq. (2). This function is the (multivariate) composition of two functions, namely,  $\theta = \theta(\mathbf{x}, k)$  and  $\text{Huber}(\text{Soft-Threshold}(\mathbf{x}, \theta))$ . In particular, the former is continuously differentiable (i.e.,  $C^1$ ) in  $\mathbf{x}$ , since it is simply a linear combination of sample mean and sample standard deviation both of which are  $C^1$  functions. To establish the theorem, we only need to show that  $\text{Huber}(\text{Soft-Threshold}(\mathbf{x}, \theta))$  is also a  $C^1$  function in  $(\mathbf{x}, \theta)$ .

By definition,  $\text{Huber}(\text{Soft-Threshold}(\mathbf{x}, \theta))$  is defined entry-wise on  $\mathbf{x}$  as

$$\text{Huber}(\text{Soft-Threshold}(x, \theta)) = \begin{cases} \delta x - \delta\theta - \frac{1}{2}\delta, & \text{if } x > \theta + \delta; \\ \frac{1}{2}(x - \theta)^2, & \text{if } \theta \leq x \leq \theta + \delta; \\ 0, & \text{if } x < \theta. \end{cases} \quad (30)$$

From here it is easy to check that  $\text{Huber}(\text{Soft-Threshold}(x, \theta))$  is continuous in  $(x, \theta)$ . Its gradient with respect to  $(x, \theta)$  is given by

$$\frac{\partial \text{Huber}(\text{Soft-Threshold}(x, \theta))}{\partial (x, \theta)} = \begin{cases} (\delta, -\delta), & \text{if } x > \theta + \delta; \\ (x - \theta, \theta - x) & \text{if } \theta \leq x \leq \theta + \delta; \\ (0, 0), & \text{if } x < \theta, \end{cases} \quad (31)$$

which is also continuous. This concludes the proof. □

## C IMPLEMENTATION DETAILS ON SPARSE MATRIX MULTIPLICATIONS

We describe how we implement sparse matrix multiplications for Spark FFN and Attention in `gemma.cpp`. We start by focusing on a batch size of one for decoding before expanding our discussion to larger batch sizes and prefill.

With batch size of 1, both Spark FFN and Spark Attention utilize two types of sparse vector-matrix multiplication: vector-masked matrix multiplication and sparse vector-matrix multiplication (Figure 2). Given a vector  $q$  and a matrix  $w$ , vector-masked matrix multiplication multiplies  $q$  with the non-masked columns of  $w$  based on a masking vector  $m$ . Masked columns yield a zero output. Sparse vector-matrix multiplication, on the other hand, involves a vector that contains many zeroes being multiplied by a dense matrix.

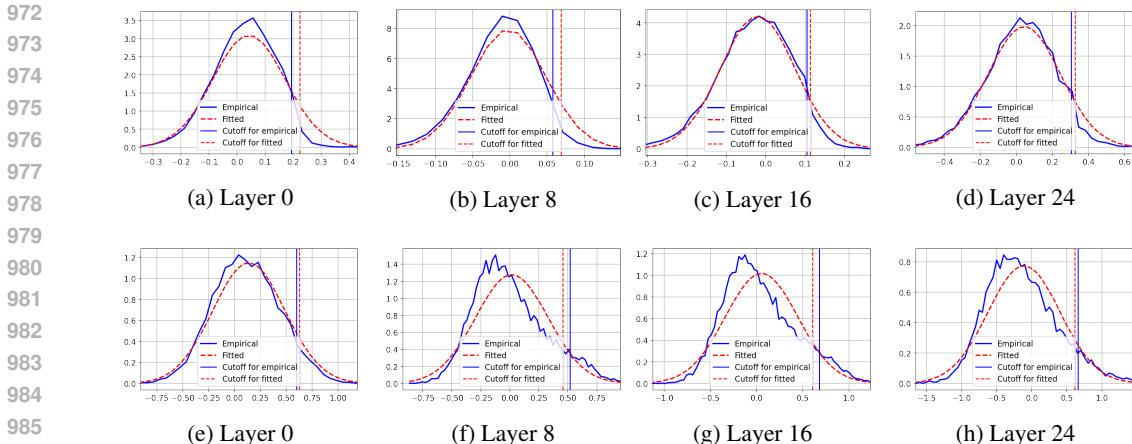


Figure 6: Distribution of the entries of the input activation to statistical top- $k$  in Spark FFN (see Figure 7 for result of Spark Attention). The two rows correspond to activation at two positions 0 and 1000 of an input, and the columns correspond to activation at four different depth levels  $\{0, 8, 16, 24\}$  of the 26-layer pretrained Spark Gemma-2. The input is the first 1000 tokens of the first essay from <https://paulgraham.com/articles.html> prepended with the BOS token. We compare the empirical distribution (*Empirical*) with the Gaussian distribution whose mean and standard deviation (std) are computed as the sample mean and std of the input (*Fitted*). We see that the Gaussian closely approximates the empirical distribution. We also compare the cutoff value estimated from the Gaussian, i.e.,  $\theta(x, k)$  used in Eq. (2) with  $k/d = 5\%$  (*Cutoff for fitted*), with the cutoff value for obtaining 8% nonzeros on the empirical distribution (*Cutoff for empirical*). It can be seen that these two cutoff values are close.

In Spark FFN, we perform vector-masked matrix multiplication for  $K^T(I - P)q$  (Figure 2a). The masking vector is generated from the output of Statistical-Top $_k(K^T Pq)$ . In the CPU implementation, the columns of  $w$  are loaded from DRAM one at a time. Based on the mask, Spark FFN skips loading the masked columns of  $w$  from DRAM and the associated computations. Spark FFN utilizes SIMD operations (as in the original Gemma implementation). To further enhance performance, Spark FFN utilizes software CPU prefetching (*builtin\_prefetch*) to overlap loading from DRAM to the CPU cache with computations.

The same masking vector also identifies the zero entries in the intermediate vector that is multiplied by matrix  $V$  (Figure 2b). For this sparse vector-matrix multiplication, we store the matrix in row format. Each CPU thread processes a tile of the matrix while skipping the loading and computation of the masked rows. Prefetching and SIMD operations are applied similarly in this context.

Spark Attention utilizes these two types of sparse matrix multiplication operators to accelerate qkv computations for each head.

When extending to decoding with batch sizes greater than one or prefill, we continue to use individual masks to skip *computations* while using a union of masks from each vector within the batch to create unified masks for *memory loading*. With larger batches, Spark transformer is expected to save less memory loading (vs. original Gemma), unless there is significant overlap in top- $k$  positions within the same batch. Nonetheless, the Spark transformer consistently reduces FLOP by skipping computations based on individual masks within the batch.

## D ADDITIONAL EXPERIMENTAL RESULTS

### D.1 DISTRIBUTION OF INPUTS TO STATISTICAL TOP- $k$

The underlying assumption for statistical top- $k$  is that the activation vector upon which it is applied to, namely, the pre-GELU activation in Spark FFN and the pre-softmax activation in Spark Attention, can be modeled as being drawn from an i.i.d. Gaussian distribution. Here we provide empirical evaluation on the distribution of these activation vectors for Spark Gemma2. Results for Spark FFN

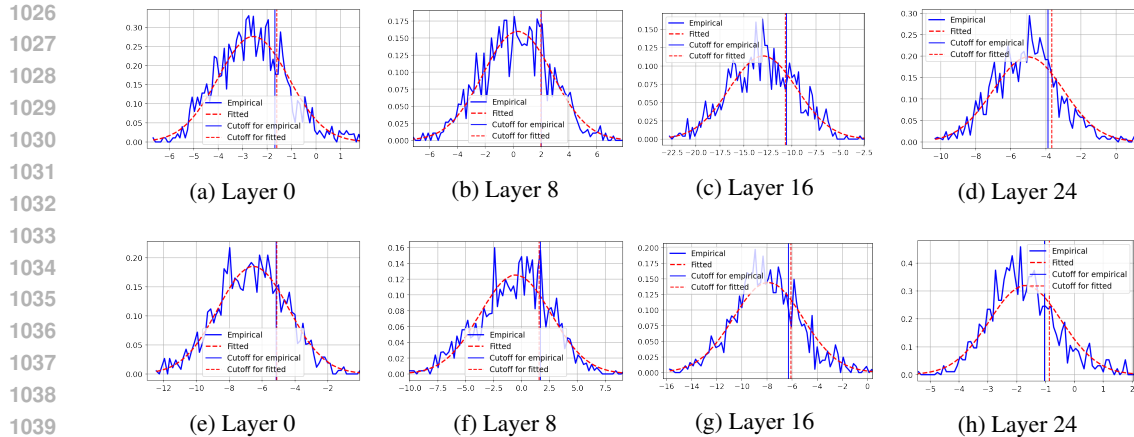


Figure 7: Distribution of the entries of the input activation to statistical top- $k$  in Spark Attention (see Figure 6 for result of Spark FFN). The two rows correspond to activation for two different attention heads, and the columns correspond to activation at four different depth levels  $\{0, 8, 16, 24\}$  of the 26-layer pretrained Spark Gemma-2. Model input is the first 1000 tokens of the first essay from <https://paulgraham.com/articles.html> prepended with the BOS token, and we examine activation of the last token (i.e., inner product between the query embedding of the 1001st token and all 1001 key embeddings). We compare the empirical distribution (*Empirical*) with the Gaussian distribution whose mean and standard deviation (std) are computed as the sample mean and std of the input (*Fitted*). We see that the Gaussian closely approximates the empirical distribution. We also compare the cutoff value estimated from the Gaussian, i.e.,  $\theta(x, k)$  used in Eq. (2) with  $k = 256$  (*Cutoff for fitted*), with the cutoff value for obtaining top 256 entries on the empirical distribution (*Cutoff for empirical*). It can be seen that these two cutoff values are close.

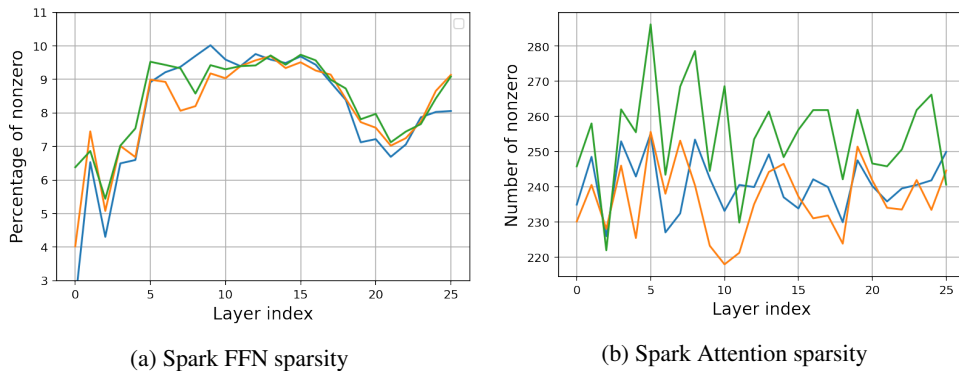


Figure 8: Sparsity in the intermediate activation of Spark FFN and Spark Attention *during evaluation* (see Figure 1 for results during training). For FFN, we use a simple prompt “test” and report the percentage of nonzero entries in generating the 5th, 10th, and 15th token. For Attention, we report the number of nonzero entries at the 512th, 1024th, and 2048th token during prefill.

and Spark Attention are provided in Figure 6 and Figure 7, respectively. The results show that the distribution holds close proximity to a Gaussian, hence justifying the use of statistical top- $k$ .

## D.2 SPARSITY LEVEL DURING EVALUATION

Complementing Figure 1 which reports sparsity level during pretraining, here we report the sparsity level during evaluation to confirm that statistical top- $k$  produces the same level of sparsity during test time. The results are presented in Figure 8 for some arbitrarily selected tokens. For Attention, in particular, we select tokens at the positions 512, 1024, and 2048 which are all above our choice of  $k = 256$  for Spark Attention.

### D.3 BATCHING ANALYSIS

Figure 9 provides the performance comparison between Spark Gemma 2 and Gemma 2, measured in prefill throughput (tokens/sec) across varying chunk sizes. We use a 4096-token prompt on a 16-core CPU VM. A similar trend is expected during the decoding phase with varying batch sizes.

Our analysis shows that Spark Gemma-2 provides the highest speedup at batch size 1, and again at large batch sizes (e.g.  $>8$ ), where the compute FLOP becomes the primary bottleneck.

For Gemma-2, as seen in the figure, increasing batch/chunk size leads to a significant improvement in prefill throughput until the batch size reaches 8. This improvement occurs because batching reduces memory access by reusing weights across multiple tokens in the CPU cache. Once the computation becomes the bottleneck (i.e. batch = 8), further batching provides diminishing returns.

In contrast, Spark Gemma-2 behaves differently. When the batch size increases from 1 to 2, we observe minimal throughput change. This is due to the lack of overlap in top-k positions between the tokens, resulting from the high sparsity. However, as the batch size increases beyond 4, Spark Gemma-2 starts benefiting from weight reuse, similar to Gemma-2. Spark Gemma-2 continues to show improvements in throughput until the batch size reaches approximately 64, where it eventually becomes FLOP-bound, much later than Gemma-2 due to the reduced FLOP requirements.

Overall, Spark Gemma demonstrates the most significant gains in two scenarios: when the batch size is 1, a common setting for desktop or mobile devices decoding, and when the batch size is large enough that FLOP becomes the dominant bottleneck.

### D.4 ADDITIONAL ABLATION STUDIES

In this section, we provide ablation studies for understanding the effect of the individual components of Spark Transformer. Towards that, we plot the training loss curves for Gemma-2 and Spark Gemma-2, see Figure 10. Here, we restrict to the first 80k training steps out of the 500k total steps since it is costly to fully train all ablation models, and that 80k steps is sufficient for seeing the trend. We can see that Spark Gemma-2 slightly lags behind Gemma-2. However, as demonstrated in Table 2, that small difference in training loss does not lead to a substantial difference in evaluation quality.

In our ablation studies below, we add a single component at a time to Gemma-2 and see the quality impact.

**Spark FFN vs Spark Attention.** To understand the effect of Spark FFN vs Spark Attention, we conduct an experiment where only attention is switched from a standard one to Spark Attention, whereas the FFN remains the standard one. The result is illustrated as *Gemma-2 + Spark Attention* in Figure 10. It can be seen that Spark Attention provides a minor quality gain over Gemma-2. In comparing *Gemma-2 + Spark Attention* with Spark Gemma-2, this also shows that further adding Spark FFN slightly hurts model quality. As noted above, such a small difference does not lead to substantial quality impact on the evaluation tasks. Hence, we conclude here that none of Spark FFN and Spark Attention has significant quality impact.

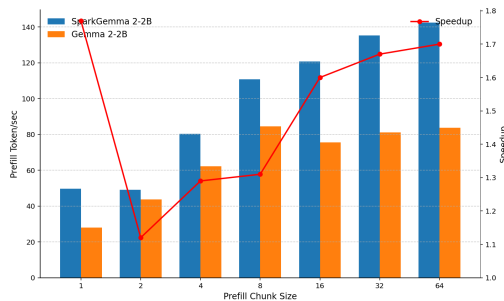


Figure 9: Spark Gemma-2 vs. Gemma-2 Prefill Token/Sec with Varying Chunk Sizes. We use a prompt length of 4096 tokens on a 16 core CPU VM.

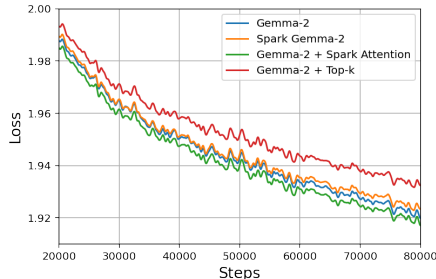


Figure 10: Ablation study in terms of training loss in the first 80k training steps (out of 500k total steps).

**Sparsity enforcing vs Low-cost predictor.** Sparsity enforcing via statistical top- $k$  and low-cost activation predictor are two relatively independent components of Spark Transformer. This means that, upon the standard Gated FFN (see Eq. (12)) that is used in Gemma-2, which we rewrite here for convenience:

$$\text{Gated-FFN}(\mathbf{q}; \mathbf{K}_1, \mathbf{K}_2, \mathbf{V}) = \mathbf{V} \cdot (\sigma(\mathbf{K}_1^\top \mathbf{q}) \odot (\mathbf{K}_2^\top \mathbf{q})), \quad (32)$$

we may choose to only apply statistical top- $k$  for enforcing sparsity, i.e.,

$$\text{Topk-Gated-FFN}(\mathbf{q}; \mathbf{K}_1, \mathbf{K}_2, \mathbf{V}) = \mathbf{V} \cdot (\sigma(\text{Statistical-Top}_k(\mathbf{K}_1^\top \mathbf{q})) \odot (\mathbf{K}_2^\top \mathbf{q})). \quad (33)$$

Note that applying a sparsifying function on the input to the nonlinear function  $\sigma(\cdot)$  as in Eq. (33) is a common choice in the literature of sparse activations, e.g., Mirzadeh et al. (2023); Song et al. (2024a); Lee et al. (2024b); the main difference between these works lies in the specific sparsity enforcing technique, see Table 4 for a summary. In addition to the sparsifying function, Spark FFN also has another architectural change for the purpose of introducing a low-cost predictor. Here, we rewrite Spark FFN for ease of comparison with Eq. (33):

$$\text{Spark-FFN}(\mathbf{q}; \mathbf{K}, \mathbf{V}, k, r) \stackrel{\text{def}}{=} \mathbf{V} \cdot (\sigma(\text{Statistical-Top}_k(\mathbf{K}^\top \mathbf{P}\mathbf{q})) \odot (\mathbf{K}^\top (\mathbf{I} - \mathbf{P})\mathbf{q})). \quad (34)$$

Analogous to FFN, we may also only add statistical top- $k$  to attention without the low-cost predictor, i.e., by switching from standard Attention in Eq. (13) to the following:

$$\text{Topk-Attention}(\mathbf{q}; \mathbf{K}, \mathbf{V}) \stackrel{\text{def}}{=} \mathbf{V} \cdot \text{softmax}(\text{Statistical-Top}_k(\mathbf{K}^\top \mathbf{q})). \quad (35)$$

Here, we aim to understand the effect of introducing statistical top- $k$  without the low-cost predictor. Towards that, we conduct an experiment where FFN and Attention in Gemma-2 are replaced with Eq. (33) and Eq. (35), respectively. The result is illustrated as *Gemma-2 + Top-k* in Figure 10. It can be seen that the training loss becomes notably larger and the gap compared to Gemma-2 is further increasing with more training steps. This result demonstrates that while the low-cost predictor is introduced to Spark Transformer for reducing the cost in predicting the nonzero entries, it also helps in bridging the gap from the introduction of statistical top- $k$ . In other words, Transformer with low-rank predictors in FFN and Attention is more amenable to activation sparsification without quality loss.

## E COMPARISON WITH RELATED WORK ON ACTIVATION SPARSITY IN FFN

In Table 4, we provide a summary of recent work on enabling sparse activation in the latest LLMs.

We can see that our Spark Transformer leads to a FLOPs reduction of -72%, which is more than all the other methods. This comes at a cost of -0.9% quality loss, which is lower than most of the other methods (i.e., ReLUification and ProSparse) and is on par with the best alternative, i.e., HiRE.

## F ADDITIONAL DISCUSSION ON STATISTICAL TOP- $k$

### F.1 NOVELTY UPON EXISTING WORK

We note that ideas similar to our statistical top- $k$  have appeared in the literature. In particular, Shi et al. (2019) introduced the idea of fitting a Gaussian distribution to the entries of an input vector

<sup>2</sup>Total training cost relative to the base model. For finetuning based approaches, such as ReLUification (on Falcon and Llama) and ProSparse, the total training cost includes both the pretraining cost and the finetuning cost.

<sup>3</sup>Quality loss relative to the base model. Here the numbers are based on the results reported in their respective papers. Note that a different set of evaluation benchmarks is used in each paper. For ReLUification, this set contains ARC-Easy, HellaSwag, Lambada (for OPT) and Arc-E, Arc-C, HellaSwag, BoolQ, PIQA, LAMBADA, TriviaQA, WinoGrande, SciQ (for Falcon and Llama). For ProSparse, this set contains HumanEval, MBPP, PIQA, SIQA, HellaSwag, WinoGrande, COPA, BoolQ, LAMBADA, and TyDiQA. For HiRE, this set contains WMT14/WMT22, SuperGLUE, Multiple QA datasets, and multiple discriminative tasks datasets. For CATS, this set contains OpenBookQA, ARC Easy, Winogrande, HellaSwag, ARC Challenge, PIQA, BoolQ, and SCI-Q. For Spark Transformer, the datasets are those reported in Table 2.

<sup>4</sup>Results reported here are for the stage 1 training of their paper.

<sup>5</sup>Specifically, -80% on odd layers only, and -60% on average.

Table 4: Comparison with related work on enforcing activation sparsity in FFN of LLMs. Spark Transformer has the largest FLOPs reduction with one of the smallest quality loss.

|   | Main Techniques                     |                         | Main Results            |                            |                   |                      |                   |
|---|-------------------------------------|-------------------------|-------------------------|----------------------------|-------------------|----------------------|-------------------|
|   | Enforce sparsity                    | Predict support         | Base model              | Training cost <sup>2</sup> | Sparsity (%zeros) | Quality <sup>3</sup> | FFN FLOPs         |
| ReLUification <sup>4</sup><br>(Mirzadeh et al., 2023) | ReLU                                | None                    | OPT 1.3B                | +0%                        | 93%               | -2%                  | -62%              |
|   | ReLU                                | None                    | Falcon 7B<br>Llama 7B   | +2%<br>+3%                 | 94%<br>62%        | -2.5%<br>-1.9%       | -62%<br>-42%      |
| ProSparse<br>(Song et al., 2024a)                     | ReLU +<br>$\ \cdot\ _1$             | None                    | Llama2 7B               | +1.8%                      | 88%               | -1.1%                | -59%              |
|   |                                     |                         | Llama2 13B              | +6.7%                      | 88%               | -1.4%                | -59%              |
|   |                                     | 2-layer FFN             | Llama2 7B<br>Llama2 13B | NA<br>NA                   | 75%<br>78%        | NA<br>NA             | NA<br>NA          |
| HiRE<br>(Yerram et al., 2024)                         | Group top <sub>k</sub> + commonpath | Low-rank / quantization | PALM2 1B                | +0%                        | 80%               | -0.8%                | -60% <sup>5</sup> |
| CATS<br>(Lee et al., 2024b)                           | Thresholding                        | None                    | Mistral 7B              | +0%                        | 50%               | -1.5%                | -33%              |
|   |                                     |                         | Llama2 7B               |                            | 50%               | -2.4%                | -33%              |
| Spark Transformer                                     | Statistical-top <sub>k</sub>        | Partial dimensions      | Gemma 2B                | +0%                        | 92%               | -0.9%                | -72%              |

and estimating a threshold from quantile functions. Then, M Abdelmoniem et al. (2021) extended the approach to additional distributions. In both cases, the method is used for solving the problem of distributed training. Here, we summarize our contribution upon these works:

- We are the first to use statistical top- $k$  for enforcing activation sparsity in Transformers. Improving Transformer efficiency via activation sparsity has become a very popular research topic (see Section 5), but may have been suffering from a lack of efficient top- $k$  algorithms for enforcing sparsity. Hence, the introduction of statistical top- $k$  may facilitate the development of this area.
- Synergizing statistical top- $k$  into Transformers is nontrivial. Since the method of statistical top- $k$  is based on fitting a statistical distribution to the activation vector, there is the need to understand the distributions of different activations in order to determine which particular activation vector is suited for the application of statistical top- $k$  and the associated choice of the distribution. In our case, we decide that statistical top- $k$  should be applied to the activation before the nonlinear function (for FFN) and before softmax (for Attention) since entries of this vector provably follows a Gaussian distribution at random initialization. We also verify empirically that statistical top- $k$  is still reliable even after initialization.
- We extended statistical top- $k$  from using the hard-thresholding operator with the estimated statistical threshold to the soft-thresholding operator. This leads to a continuous optimization landscape that may have facilitated the optimization. Empirically, we found this choice to provide quality benefits for Spark Transformer.
- We provide the first theoretical justification for the correctness of statistical top- $k$ , see Theorem 1.
- We reveal the conceptual connection between statistical top- $k$  and several related top- $k$  operators in the literature, see Section 2.3. Such connections may motivate the development of more powerful top- $k$  algorithms in the future.

## F.2 HANDLING CASES WHEN THE ACTIVATION IS SHARDED

The training of modern large Transformer models usually requires sharding certain model weights and activations across multiple computation devices, due to physical limitations on each device’s memory. In particular, if sharding is used for the activations upon which the statistical top- $k$  is applied to, i.e., the pre-GELU activation in FFN and the pre-softmax activation in attention, extra care needs to be taken so that statistical top- $k$  is applied correctly. While this has not been the case



1242 for our experiment on Gemma-2 2B, here we discuss possible solutions if this case arises, e.g., when  
1243 training a larger Spark Transformer for which sharding relevant activations may become necessary.

1244 Assume that an activation vector of length  $N$  is sharded over  $m$  devices, and we are interested in  
1245 finding approximately the top- $k$  entries of  $N$  with the largest value. There are two ways of applying  
1246 statistical top- $k$  for this purpose.

- 1247
- 1248 • Global statistical top- $k$ . Here we require each device to compute the mean and variance for entries  
1249 on itself, then communicate them to all other devices. In this case, each device receives  $m - 1$   
1250 mean and variance values, which can be combined with mean and variance on its own to obtain  
1251 global mean and global variance. Then, the rest of the steps in statistical top- $k$  can be carried out  
1252 on individual devices. In this method, the output is exactly the same as if applying statistical top- $k$   
1253 without activation sharding. In terms of cost, there is extra computation coming from aggregating  
1254 mean and variance from all devices, but the cost is very low as it requires only  $O(k)$  FLOPs.  
1255 The method also introduces a communication cost, but the cost is again small as each device only  
1256 needs to send / receive  $2k - 2$  floating point numbers.
- 1257 • Local statistical top- $k$ . Here we simply apply statistical top- $k'$  to entries on its own device with  
1258  $k' = k/m$ . The method is sub-optimal in the sense that it does not necessarily produce the same  
1259 output as applying the global statistical top- $k$ . However, it has the benefit of not adding any  
1260 computation and communication cost.

1261 In cases where  $k \ll N$ , the global statistical top- $k$  above is cheap enough and hence could be a  
1262 natural choice.

1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295