EFFICIENT LLM INFERENCE USING DYNAMIC INPUT PRUNING AND CACHE-AWARE MASKING

Marco Federici* ¹ Davide Belli* ¹ Mart van Baalen ¹ Amir Jalalirad ¹ Andrii Skliar ¹ Bence Major ¹ Markus Nagel ¹ Paul Whatmough ¹

ABSTRACT

While mobile devices provide ever more compute power, improvements in DRAM bandwidth are much slower. This is unfortunate for large language model (LLM) token generation, which is heavily memory-bound. Previous work has proposed to leverage natural dynamic activation sparsity in ReLU-activated LLMs to reduce effective DRAM bandwidth per token. However, more recent LLMs use SwiGLU instead of ReLU, which result in little inherent sparsity. While SwiGLU activations can be pruned based on magnitude, the resulting sparsity patterns are difficult to predict, rendering previous approaches ineffective. To circumvent this issue, our work introduces Dynamic Input Pruning (DIP): a predictor-free dynamic sparsification approach, which preserves accuracy with minimal fine-tuning. DIP can further use lightweight LoRA adapters to regain some performance lost during sparsification. Lastly, we describe a novel cache-aware masking strategy, which considers the cache state and activation magnitude to further increase cache hit rate, improving LLM token rate on mobile devices. DIP outperforms other methods in terms of accuracy, memory and throughput trade-offs across simulated hardware settings. On Phi-3-Medium, DIP achieves a 46% reduction in memory and 40% increase in throughput with < 0.1 loss in perplexity. We plan to open source the code for HW simulator, methods, and experiments in this paper.

1 Introduction

Specialized neural processing unit (NPU) accelerators for AI have become standard in recent smartphone SoCs. These processors enable efficient neural network inference ondevice, opening the way for powerful on-device AI applications. However, as shown in Figure 2, memory size has not improved at the same pace as processing power (Jhang et al., 2021; Kwon et al., 2023; Gholami et al., 2024). This memory bottleneck is especially evident for LLMs, which grew to billions and trillions of parameters, and for which all weights must be accessed to generate each individual output token. For instance, a model with 14B parameters like Phi-3-Medium quantized at INT4 precision occupies approximately 7GB, which approaches the average smartphone *total* DRAM capacity (Counterpoint, 2021; Grossi & Lorenz, 2024).

ArXiv preprint, currently under review.

Snapdragon branded products are products of Qualcomm Technologies, Inc. and/or its subsidiaries.

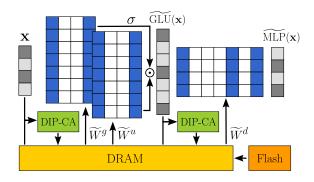


Figure 1. Overview of the Dynamic Input Pruning method in a Gated MLP. The Cache-Aware masking component outputs a column sparsity mask based on current input activations and weights currently cached in DRAM. The unpruned columns are retrieved from DRAM cache, if available, or loaded from Flash in case of cache misses.

But, in fact, the OS and active applications already a significant portion of the total DRAM capacity, typically leaving a few GBs for an on-device LLM.

To relax the DRAM capacity pressure, it has been previously proposed to load a model directly from flash, instead of first copying it in DRAM (Alizadeh et al., 2023; Xue et al., 2024). By leveraging *dynamic sparsity* and introducing *DRAM caching* the latency overhead of loading a

^{*}Equal contribution ¹Qualcomm AI Research. Qualcomm AI Research is an initiative of Qualcomm Technologies, Inc. Correspondence to: Marco Federici <mfederic@qti.qualcomm.com>, Davide Belli <dbelli@qti.qualcomm.com>, Mart van Baalen <mart@qti.qualcomm.com>, Markus Nagel <markusn@qti.qualcomm.com>.

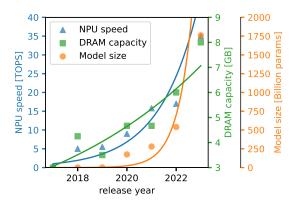


Figure 2. Comparing the largest LLMs by year of release (Zhao et al., 2023) against NPU and DRAM trends in recent iPhone devices². An exponential fit highlights the linear increase in DRAM size opposed to the exponential increases in NPU processing power and model sizes.

model from flash can be mitigated significantly. Dynamic sparsity (Mirzadeh et al., 2023; Lee et al., 2024; Akhauri et al., 2024) exploits the fact that all operations during token generation are matrix-vector multiplications, and the fact that in LLMs with ReLU activations, vectors contain many zero elements. For elements in the activation vectors that are zero, we can skip loading the corresponding weights in the following layer. If sparsity predictors such as those introduced in DejaVu (Liu et al., 2023) are used to predict which elements in the activation vectors will be zero, the weights that generate zero-valued activation elements can also be ignored in the preceding layer. Additionally, a DRAM cache is used to cache active parameters from previous tokens. When a DRAM cache is used, parameters only need to be loaded from flash in case of a cache miss, i.e., when an active parameter is not already in cache from previous tokens. These approaches allow even Mixtral8x7B to run on mobile devices at reasonable token generation rates (Xue et al., 2024).

Prior approaches rely on predictors that can anticipate the natural dynamic activation sparsity. As we will show in more detail in Section 3, SwiGLU MLP used in most modern LLMs (Jiang et al., 2023; Bai et al., 2023; Abdin et al., 2024; Dubey et al., 2024) have little to no natural activation sparsity. While these activations can still be pruned based on magnitude, it turns out to be much harder to predict which activations should be pruned in SwiGLU LLMs than in ReLU LLMs (see Figure 6). This issue is circumvented in previous works (Mirzadeh et al., 2023; Song et al., 2024a;b) by replacing SiLU activations with ReLU activations, and fine-tuning the LLM to adapt to the architecture change. However, this fine-tuning step is expensive, often requiring

billions of tokens of fine-tuning data to achieve reasonable accuracy. In this work we introduce Dynamic Input Pruning as a predictor-free alternative to dynamic sparsity that does not rely on expensive re-training. We also show that small LoRA adapters (Hu et al., 2021) can be used to mitigate the performance loss due to dynamic pruning.

The high natural sparsity rates of ReLU LLMs (up to 90%) have two implications: 1) for each token, only a few active parameters need to be loaded from Flash, and 2) active parameters from multiple previous tokens can be cached, leading to higher cache hit rate and, as a result, reduced flash access. SwiGLU LLMs, however, can be pruned dynamically only up to approximately 50% before model accuracy is significantly impacted (see Section 3). Hence, to reach reasonable token generation rates at lower sparsity levels, we introduce a *cache-aware* variant of dynamic input pruning. This variant takes into account the current state of the cache when deciding which parameters to use. By biasing the parameter selection towards the ones that are already in cache, we can increase cache-hit rate and hence LLM throughput, while minimally affecting model accuracy.

The main contributions in this paper are three-fold:

- 1. We analyze the limitation of methods like DejaVu (Liu et al., 2023) and "LLM in a Flash" (Alizadeh et al., 2023) in modeling activation sparsity for modern non-ReLU LLMs, demonstrating the impracticality of training predictors for SwiGLU-based models (Section 3).
- 2. We propose a predictor-free method to dynamically prune MLP layers which outperforms previous static and dynamic pruning methods on the latest LLMs across various sparsity levels. Additionally, we show that a lightweight LoRA component can further reduce the accuracy lost with sparsification. (Section 4).
- 3. We implement a cache-aware masking strategy which modifies the sparsity mask by prioritizing cached weights, achieving better trade-offs between model accuracy and latency at different DRAM constraints and hardware settings (Section 5).

2 BACKGROUND AND RELATED WORK

This section first reviews quantization and static sparsity, two alternative approaches to reduce model size, and later discusses prior works employing dynamic sparsity.

Static Pruning and Quantization Pruning and quantization techniques have been extensively explored to minimize the memory footprint of neural networks (Kuzmin et al., 2023). Unstructured pruning typically involves removing individual weight parameters, whereas structured

²Data collected from Wikipedia (2024a). Values are averaged in case of multiple model releases per year.

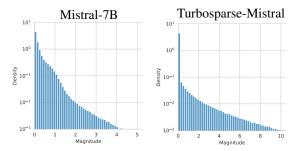


Figure 3. GLU Activation Magnitude distribution for Layer 31 Mistral-7B (left), and the corresponding re-trained ReLU-fied Turbosparse Mistral (right).

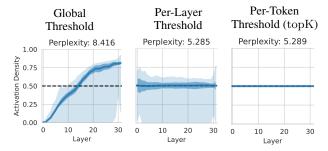


Figure 4. Layer activation density for different GLU thresholding strategies at 50% MLP target density on a Mistral-7B model. Darker shaded region indicate the standard deviation, while lighter shading covers minimum and maximum values. Setting a threshold for each layer or each token result in similar perplexity. The dense model perplexity is 5.249.

pruning eliminates entire weight structures (van der Ouderaa et al., 2023; Frantar & Alistarh, 2023; Sun et al., 2023). Alternative approximation strategies involve low bit-width quantization (Nagel et al., 2021; Frantar et al., 2022; van Baalen et al., 2024; Huang et al., 2024) and low-rank representations (Zhang & Papyan, 2024). All prior techniques, however, tend to degrade performance, which can be partially mitigated through model fine-tuning. This fine-tuning process necessitates a representative dataset and significant computational resources, and it is therefore prohibitive for certain applications. Another downside of static pruning is its limited flexibility, as the set of retained weights is fixed for all activations, while the optimal set of weights may vary across different tokens.

Dynamic Sparsity Recent work proposed dynamic sparsity methods (Mirzadeh et al., 2023; Liu et al., 2023; Lee et al., 2024; Akhauri et al., 2024) to select active weights based on the current inputs, and thus overcoming part of the limitations in static methods.

Recent approaches (Mirzadeh et al., 2023; Alizadeh et al., 2023; Song et al., 2024b) aim to reduce the number of MLP parameters loaded at inference time by exploiting the sparse activation pattern of the intermediate ReLU activations. Liu

et al. (2023) (*DejaVu*) further showed that, in LLM with ReLU activations such as OPT (Zhang et al., 2022) and Llama v1 (Roumeliotis et al., 2023), even small predictors can accurately determine active units with moderate training and memory overhead. Whenever the predictions are correct, the DejaVu approach does not introduce any approximation error since only rows and column corresponding to zero entries are removed.

Other recent work extended the approach from DejaVu aiming to reduce system latency and memory overhead. Akhauri et al. (2024) replaced the individual predictors with a single, global predictor to determine the active units. Alizadeh et al. (2023) added an LFU cache on the MLP parameters and introduced weight bundling to reduce memory overhead. Song et al. (2023) proposed identifying most active neurons offline and pre-loading them on GPU for faster access, while loading the remaining neurons on-demand based on predictor outputs. Xue et al. (2024) builds an engine to optimize LLM inference given certain input hardware specifications. Interestingly, this work shows that neuron pre-loading is not necessary to achieve high system throughput.

Most dynamic sparsity approaches inherently rely on the sparse activation pattern induced by the ReLU non-linearity, however, the vast majority of recent LLM architectures (Jiang et al., 2023; Abdin et al., 2024; Dubey et al., 2024) use non-linearities that do not naturally induce many hard zeros (cf. Figure 3). While *ReLU-fication*, the process of replacing the original network activations with ReLU activations, has shown promising results (Mirzadeh et al., 2023; Song et al., 2024b), though on the other hand it requires a computationally demanding retraining procedure which may fail to replicate the original performance.

For this reason, in the following sections we investigate the effectiveness of dynamic sparsity strategies when applied directly to modern LLMs.

3 DYNAMIC SPARSITY ON MODERN LLMS

Modern LLMs consist of alternating sequences of *attention* and *MLP* blocks. In most cases, the majority of network weights resides in the latter component. This became even more pronounced with the introduction of *Group Query Attention* (Roumeliotis et al., 2023), reducing the number of parameters in attention layers. MLP architectures in recent LLMs (Jiang et al., 2023; Abdin et al., 2024; Dubey et al., 2024) frequently consist of a *Gated non-Linear Unit* (GLU) layer, which maps the activations into high-dimensional feature space followed by a linear down projection:

$$MLP(\mathbf{x}) = \mathbf{W}^d GLU(\mathbf{x}) \tag{1}$$

$$GLU(\mathbf{x}) = \mathbf{W}^u \mathbf{x} \odot \sigma(\mathbf{W}^g \mathbf{x}) \tag{2}$$

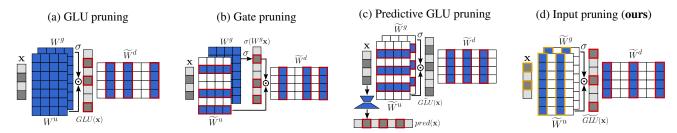


Figure 5. Diagrams describing the sparsity schemes for GLU pruning, Gate pruning, Predictive pruning and Input pruning. Activations with larger magnitude are visualized in darker colors, while dark blue indicates active weights for the current input.

in which W^u , W^g , W^d will be referenced as up, gate, and down weight matrices, respectively. The non-linearity σ often consists of a SiLU function (Dauphin et al., 2017; Elfwing et al., 2018; Naveed et al., 2023).

We define as follows the sparse version of an MLP:

$$\widetilde{\mathrm{MLP}}(\mathbf{x}) = \widetilde{\boldsymbol{W}}^d \left(\widetilde{\boldsymbol{W}}^u \mathbf{x} \odot \sigma(\widetilde{\boldsymbol{W}}^g \mathbf{x}) \right), \tag{3}$$

where the tilde indicates that the original weight matrices are replaced with a sparse version, which may result in an approximation error.

As discussed in Section 2, SwiGLU MLPs exhibit no natural activation sparsity, and existing methods to circumvent this are expensive and time-consuming. In the remainder of this section we will first show that non-sparse activations can still be pruned based on magnitude. We then discuss how this can be used to dynamically prune weights. Lastly, we demonstrate that predictor-based approaches such as those introduced by DejaVu (Liu et al., 2023) and its derivatives fail for dynamic, magnitude-based sparsity.

3.1 Magnitude-based Dynamic Pruning

Since the number of hard zeros in SwiGLU architectures is negligible, to enable higher sparsity we first relax the condition of removing only zeros, and instead prune values that are smaller than a predetermined threshold t.

In Figure 4, we compare the effect of three GLU thresholding strategies on Mistral-7B layer activations in terms of WikiText-2 perplexity: a *global threshold*, i.e., a fixed threshold *t* is used to prune the activations in all layers; a *per layer* threshold, i.e. a fixed threshold which is set independently per layer based on the CDF of activations over a small calibration set, but shared over different tokens; and a *per-token* threshold, which is based on the CDF of the current activations of each layer dynamically. Note that the latter corresponds to simply keeping the top-K largest magnitude values for each token. In all cases the thresholds are chosen to ensure 50% of activations are kept on average.

We find that a global threshold performs poorly compared to a per-layer or per-token threshold. This is due to the fact that activation magnitude varies considerably per layer and thus a fixed threshold results in too aggressive pruning for initial layers. Lastly, we find that the per-token pruning strategy performs comparably to the per-layer strategy, but with the additional benefit of pruning a constant number of activations for each token. For this reason, we only consider top-K pruning in the remainder of this paper.

3.2 Strategies for Dynamic Sparsity

We now consider several dynamic weight sparsification schemes that exploit dynamic magnitude-based activation pruning. Figure 5a-c shows the different schemes for GLU LLMs described in the remainder of this section.

GLU pruning (Figure 5a) When values of GLU(x) are pruned, the corresponding columns of W^d can also be pruned as they will be multiplied with a zero:

$$\widetilde{\boldsymbol{W}}_{:,i}^{d} = \begin{cases} \mathbf{0} & \text{if } i \notin \text{topK}\left(|\text{GLU}(\mathbf{x})_{i}|\right) \\ \boldsymbol{W}_{:,i}^{d} & \text{o.w.} \end{cases}$$
 (4)

Note that, instead of top-K pruning, other thresholding methods (such as those presented in Section 3.1) can be used instead

This strategy can sparsify only one of the three large MLP matrices, reaching at most 33% MLP sparsity.

Gate pruning (Figure 5b) Assuming that most of the small entries are determined by the non-linearity σ , one could first perform a dense matrix multiplication to determine $\sigma(\mathbf{W}^g\mathbf{x})$, then sparsify both \mathbf{W}^u and \mathbf{W}^d based on the small (partial) activations:

$$\widetilde{\boldsymbol{W}}_{:,i}^{u}, \widetilde{\boldsymbol{W}}_{i}^{d}, = \begin{cases} \mathbf{0}, \mathbf{0} & \text{if } i \notin \text{topK}\left(|\sigma(\boldsymbol{W}^{g}\mathbf{x})|\right) \\ \boldsymbol{W}_{:,i}^{u}, \boldsymbol{W}_{i}^{d} & \text{o.w.} \end{cases}$$
(5)

Similarly, we can define $Up\ Pruning$ as the equivalent strategy in which the activations $W^u x$ are used to prune W^g and W^d , respectively. In contrast to GLU pruning, Gate (and Up) pruning can prune up to 66% of the MLP weights. However, this pruning strategy may be sub-optimal since

the network parameters are pruned based on partial GLU activations.

Predictive GLU pruning (Figure 5c) The same predictive pruning strategy introduced by DejaVu can be extended to GLU LLMs by training a small predictor to identify lowmagnitude GLU activations:

$$\widetilde{\boldsymbol{W}}_{:,i}^{u}, \widetilde{\boldsymbol{W}}_{:,i}^{g}, \widetilde{\boldsymbol{W}}_{i}^{d} = \begin{cases} \mathbf{0}, \mathbf{0}, \mathbf{0} & \text{if } i \notin \text{topK}(pred(\mathbf{x})) \\ \boldsymbol{W}_{:,i}^{u}, \boldsymbol{W}_{:,i}^{g}, \boldsymbol{W}_{i}^{d} & \text{o.w.} \end{cases}$$
(6)

Here $pred(\mathbf{x})$ indicates the logits of a predictor trained to identify the smallest activations using cross-entropy loss.

Predictive GLU pruning methods can sparsify all MLP weight matrices, potentially resulting in higher sparsity levels at the same model accuracy. However, as we will show next, it is extremely challenging to train predictors for GLU pruning.

3.3 The challenge of predicting GLU activations

As demonstrated in Figure 6 (blue line), GLU pruning with topK thresholding allows us to prune large portions of activations with a minor impact on downstream performance. On SwiGLU networks, such as Mistral-7B, it is possible to prune up to 50% of the activations before observing a 1% decrease in MMLU performance. However, the corresponding ReLU-fied network, TurboSparse-Mistral, is much more sparsifiable, and GLU pruning can reach more than 90% activation sparsity before observing the same performance deterioration. This showcases that SwiGLU activations are inherently more sensitive to pruning.

Without any predictive mechanism, the effectiveness of GLU pruning is further hindered by the fact that only one of three MLP matrices is sparse (Figure 5a). We therefore studied the effect of predictive GLU pruning, which employs predictors trained with cross-entropy loss to identify small GLU activations. Consistently with Liu et al. (2023), after training we prune the GLU activations corresponding to the smallest outputs of the predictor $pred(\mathbf{x})$ on a per-token basis. The results are visualized in Figure 6 (orange line), which compares the effectiveness of predictive pruning on SwiGLU and ReLU-fied counterparts for the same set of predictors and training procedure.

MMLU accuracy at 50% activation sparsity drops by almost 10% on Mistral-7B, while the performance drop is negligible on the ReLU-fied counterpart. The wider gap between GLU pruning and GLU predictive pruning hints that predicting activation magnitude is a much harder task in SwiGLU networks. We hypothesize that this phenomenon is due to the interplay of various aspects. First, since SwiGLU activations are less sparse, predictor mistakes are more likely to prune relevant activations. Secondly, the activation pat-

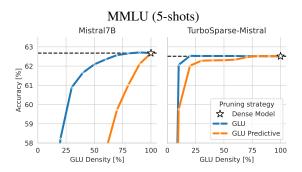


Figure 6. Comparison of the effectiveness of GLU sparsification methods on Mistral-7B (left) and its ReLU-fied counterpart (right). The predictor training procedure, architecture, and training data is equivalent across the two LLMs.

tern in GLU MLPs is less predictable than the one in the ReLU-fied counterparts. We conjecture that this is because predicting ReLU zeros is essentially sign prediction for a linear projection, a relatively easy task, whereas predicting magnitude is a more challenging task due to the interplay between the gating components.

Due to the inherent limitations of GLU pruning and the challenges of predicting complex SwiGLU activation patterns with small predictors, we consider a simple alternative approach that aims to sparsify all three MLP matrices without requiring any additional predictor architecture.

DYNAMIC INPUT PRUNING

In this section we introduce the Dynamic Input Pruning (**DIP**) strategy, which addresses the limitations of existing MLP dynamic pruning strategies on SwiGLU networks. The DIP method, visualized in Figure 5d, hinges on removing small activations in both the input of each MLP block and the intermediate GLU block using per-token topK thresholding, and thus requires no separate predictors. DIP is defined as:

$$\widetilde{\boldsymbol{W}}_{:,i}^{u}, \widetilde{\boldsymbol{W}}_{:,i}^{g} = \begin{cases} \mathbf{0}, \mathbf{0} & \text{if } i \notin \text{topK}(|\mathbf{x}|) \\ \boldsymbol{W}_{:,i}^{u}, \boldsymbol{W}_{:,i}^{g} & \text{o.w.} \end{cases}$$

$$\widetilde{\boldsymbol{W}}_{:,i}^{d} = \begin{cases} \mathbf{0} & \text{if } i \notin \text{topK}\left(|\widetilde{\text{GLU}}(\mathbf{x})|\right) \\ \boldsymbol{W}_{:,i}^{d} & \text{o.w.} \end{cases}$$
(8)

$$\widetilde{\boldsymbol{W}}_{:,i}^{d} = \begin{cases} \mathbf{0} & \text{if } i \notin \text{topK} \left(|\widetilde{\text{GLU}}(\mathbf{x})| \right) \\ \boldsymbol{W}_{:,i}^{d} & \text{o.w.} \end{cases}$$
(8)

GLU(x) corresponds to the GLU activations obtained using $\widetilde{\boldsymbol{W}}^u$ and $\widetilde{\boldsymbol{W}}^g$ instead of \boldsymbol{W}^u and \boldsymbol{W}^g . Note that we can use different levels of sparsity for the up, gate and down layers. The optimal proportion of active units is empirically determined as a function of the overall target MLP sparsity (see Appendix B.1 for details).

When compared to Predictive GLU pruning, Dynamic Input Pruning trades off the error caused by the predictor for

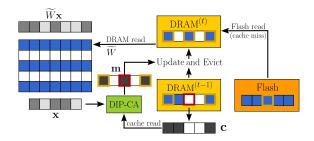


Figure 7. Diagram for the Cache-aware Dynamic Input Pruning method applied to an individual linear layer. DIP-CA uses cache state and input activations to output a sparsity mask. Flash reads are performed to load necessary weights missing from cache (marked in red), and unused weights are evicted if the DRAM is full. The sparse weight matrix is then used to process the current activations.

approximation error introduced by the use of an approximate gating layer.

Light-weight LoRA fine-tuning To compensate for the approximation error caused by the subsequent pruning of multiple MLP activations, similarly to contemporary literature (Ma et al., 2023; Sun et al., 2023; Lee et al., 2024), we add fine-tuned LoRA adapters to the sparsified up, gate and down MLP matrices. To ensure that the LoRA adapters do not lead to any additional overhead, each adapter is applied before column selection:

$$\underbrace{\boldsymbol{W}_{:,i}}_{\text{original column}} \rightarrow \underbrace{(\boldsymbol{W} + \boldsymbol{A}\boldsymbol{B})_{:,i}}_{\text{LoRA adapted column}}, \tag{9}$$

in which \boldsymbol{A} and \boldsymbol{B} refer to low-rank learnable matrices. This allows us to fuse the adapters to the original matrices after fine-tuning to avoid any memory overhead.

5 IMPROVED CACHING IN DYNAMIC SPARSITY

An effective dynamic sparsity method is not enough to guarantee good LLM inference latency on a memory-constrained device. The number of Flash reads, which bottlenecks the overall throughput, largely depends on the cache eviction policy employed to dynamically load and unload weights in DRAM during inference. In this section we discuss existing cache eviction policies and present a novel cache-aware masking strategy for dynamic sparsity. Our approach allows us to reach better accuracy-latency trade-offs with respect to solely regulating the sparsity level.

5.1 Cache eviction policies

A cache management system can be employed to optimize the decision of which parameters to evict, when the cache is full and we want to allocate new weights. Commonly used cache eviction policies include Least Recently Used (LRU) and Least Frequently Used (LFU). In LRU, the item that was last used farthest in the past is evicted. In LFU, a counter is kept (for the current session or across sessions) to track usage frequency and evict the most rarely used item.

An optimal cache eviction policy can be implemented if future DRAM reads were known in advance. Belady (1966) describes a greedy oracle algorithm to maximize the cache hit-rate over a known, fixed sequence of memory reads. Based on this optimal algorithm, previous work (Lykouris & Vassilvitskii, 2021; Liu et al., 2022; Sadek & Elias, 2024) explores methods to use estimators (Jain & Lin, 2016; Shah et al., 2022) or learnable ML models based on cache history to predict the future cache usage. While these methods outperforms simpler caching eviction policies, they require an additional training step, and their performance is upperbounded by the optimal Belady caching algorithm.

Caching methods for dynamic sparsity have also been employed, with the eviction policy used to determine which parameters should be kept in DRAM to increase the cache hit rate and therefore minimize latency. For example, Song et al. (2023) use a static cache by preloading the most active neurons, Alizadeh et al. (2023) introduce an LRU eviction policy to prune unused weights, and Xue et al. (2024) propose neuron bundling over multiple weight matrices.

5.2 Cache-aware neuron sparsity

Differently from other caching applications, the current choice of sparsity mask is not a hard constraint. We argue that suboptimal masks can be preferred to increase the system throughput while sacrificing some model accuracy. We thus propose a *cache-aware* sparsification method to strike a better trade-off between accuracy and latency. This is implemented as a re-weighting mechanism before selecting active neurons, which introduces a penalty multiplier term on the neurons activation magnitude:

$$\mathbf{s}^{(t)} = \frac{\left|\mathbf{x}^{(t)}\right| \left(\mathbf{c}^{(t-1)} + \gamma(1 - \mathbf{c}^{(t-1)})\right)}{\left\|\mathbf{x}^{(t)}\right\|_{\infty}}, \quad (10)$$

where $\mathbf{x}^{(t)}$ refers to the current input activations, $\mathbf{c}^{(t-1)}$ is a binary mask representing which neurons are currently loaded in cache, and γ is a hyper-parameter to penalize the selection of non-cached neurons. The resulting scores $\mathbf{s}^{(t)}$ are used instead of the original magnitude for dynamic pruning as discussed in Section 4. The normalization term in the denominator ensures that the re-weighting mechanisms is not sensitive to changes in the dynamic range of the activations over different tokens.

Our implementation choice is motivated by the observation that for each layer and token, few neurons fire very strongly, while activations for most neurons are orders of magnitude lower (see Figure 10, left), and have less impact on the layer **Algorithm 1** Dynamic Input Pruning with Cache-Aware re-weighting for a single linear layer.

$$\begin{aligned} & \mathbf{function} & \text{ DIP-CA } \left(\mathbf{x}^{(t)}, \mathbf{c}^{(t-1)}, \mathbf{W}, \gamma, k \right) : \\ & \mathbf{s}^{(t)} = \frac{\left| \mathbf{x}^{(t)} \right| \left(\mathbf{c}^{(t-1)} + \gamma (1 - \mathbf{c}^{(t-1)}) \right)}{\left\| \mathbf{x}^{(t)} \right\|_{\infty}} & \triangleright \textit{Eq.10} \\ & \mathbf{m}^{(t)} = \text{topK } \left(\mathbf{s}^{(t)}, k \right) \\ & \mathbf{c}^{(t)} = \text{cache_update } \left(\mathbf{c}^{(t-1)}, \mathbf{m}^{(t)} \right) \\ & \widetilde{\boldsymbol{W}}^{(t)} = \boldsymbol{W} \odot \mathbf{m}^{(t)} \\ & \mathbf{return} & \widetilde{\boldsymbol{W}}^{(t)}, \mathbf{c}^{(t)} \end{aligned}$$

outputs. The penalty term is thus set high enough to allow reranking of the activations in the intermediate regime, but low enough to not impact the strongest activations even if not in cache. We empirically find that a simple LFU eviction policy with cache-aware masking results in better latency-accuracy trade-offs even compared against the oracle algorithm from Belady (1966) without re-weighting (see Section 6.4).

We name **DIP-CA** the Cache-Aware variant of Dynamic Input Pruning, and show it outperforms DIP in terms of throughput at all perplexity levels. Figure 7 visualizes the DIP-CA procedure, which is described in detail in Algorithm 1 for a single linear layer.

6 EVALUATION

6.1 Experimental setup

A three-way trade-off in KPIs We observe how methods for efficient LLM inference must be evaluated based on three different metrics: 1) model accuracy, 2) memory footprint and 3) inference throughput. Depending on the target application, a different trade-off between these metrics can be considered optimal. To simplify the experimental analysis, we evaluate at most two metrics at a time, while ignoring or fixing the remaining ones. We first compare model performance at a fixed sparsity rate, that is at the same constraint in terms of memory footprint. The model accuracy is reported in terms of token perplexity or 5-shot accuracy on downstream tasks. We then evaluate the throughput achievable at different operating points in terms of both memory footprint and loss in model accuracy. We include in Appendix ablations on hardware specifications to show how different methods perform in various use-cases.

HW Simulator In order to accelerate the evaluation procedure, we implement a software tool to simulate the data transfers between Flash, DRAM and processing units, and then compute the latency incurred in these processes. We also implement simulators for different cache eviction poli-

cies on DRAM. This allows us to efficiently estimate the impact of different dynamic sparsity solutions on system latency, without having to implement or modify a full system on-device. As inference in the token generation scenario is bound by DRAM and Flash accesses (Jhang et al., 2021; Kwon et al., 2023; Gholami et al., 2024), we do not simulate NPU inference times, which can be largely parallelized or are negligible when compared to the more expensive memory operations. In Appendix A we include a detailed description of the simulator, including motivation for the selected parameters.

Unless otherwise stated, results throughout the paper are based on simulations for Apple A18 with DRAM I/O speed of 60 GB/s, and Flash read speed of 1 GB/s. We also conduct an ablation study on different choices for DRAM size and Flash read speed. Results are reported in Appendix D.

Methods and Models We evaluate the proposed approach against the Gate Pruning, Up Pruning, and GLU Pruning methods presented in Section 2. We further consider CATS (Lee et al., 2024) and DejaVu (Liu et al., 2023): two recent dynamic sparsity solutions. Following the original approach, each predictor consist of a neural network of 1000 hidden units, which is trained for at most 20 epochs on activations collected from a calibration set using cross entropy loss. For SwiGLU LLMs, the binary targets are generated by selecting the 10% largest activations for each token. We further compare Dynamic Input Pruning with unstructured and semi-structured (2:4, 4:8) SparseGPT pruning (Frantar & Alistarh, 2023) and quantization methods GPTQ (Frantar et al., 2022) and GPTVQ (van Baalen et al., 2024).

In DIP, we empirically calibrate the proportion of active units for a target MLP sparsity as described in Appendix B.1. For DIP-CA, we tune the cache-aware re-weighting hyperparameter on WikiText-2 validation set and fix it to $\gamma=0.2$ for all experiments. We include an ablation showing the impact of changes in γ . Following previous work (Hu et al., 2021; Lee et al., 2024), we train LoRA adapters with rank 32 for a total of 1000 iterations using a knowledge distillation loss to match the original (dense) model logits. The training procedure is consistent for both DIP and CATS models. The adapters are applied to all the up, down and gate matrices for DIP, while only up and down matrices are adapted for the CATS method.

We evaluate each method on four recent SwiGLU LLMs with different dimensionalities: Phi-3-Medium and Phi-3-Mini (Abdin et al., 2024), Llama 3 8B (Dubey et al., 2024), Mistral 7B (Jiang et al., 2023).

Datasets We use the SlimPajama (Soboleva et al., 2023) dataset as calibration datasets to train predictors and finetune LoRA adapters. We evaluate all methods on token per-

	WikiText-2 (Perplexity ↓)			MMLU (5-shot accuracy ↑)				
	Phi3Med	Phi3Mini	Llama8B	Mistral7B	Phi3Med	Phi3Mini	Llama8B	Mistral7B
Dense	4.29	6.01	6.14	5.25	78.14	70.62	65.30	62.68
GLU Pruning (oracle)	4.45	6.10	6.30	5.29	77.98	70.24	64.77	62.09
SparseGPT (unstructured)*	5.67	7.44	7.71	5.81	73.97	62.15	57.72	58.16
SparseGPT (2:4)*	6.87	9.84	10.05	7.31	68.07	53.39	43.69	51.15
SparseGPT (4:8)*	6.33	8.43	8.64	6.41	70.34	57.98	50.78	54.64
Gate Pruning	11.28	14.83	29.26	96.37	66.08	47.90	27.23	40.39
Up Pruning	7.41	10.94	13.93	7.13	70.68	56.82	46.34	50.09
DejaVu*	6.15	8.76	9.43	6.24	68.99	54.86	49.02	54.34
CATS	8.34	10.83	28.24	60.48	71.08	55.64	27.34	38.96
CATS+LoRA	5.54	7.44	16.53	9.18	72.23	60.16	30.52	44.84
DIP	5.52	7.04	7.38	5.61	75.54	66.51	60.38	59.66

7.26

5.57

75.89

Table 1. Experimental results for dynamic sparsity methods at 50% MLP sparsity. *The MLP density estimation ignores the overhead introduced by the predictors or additional data-structures. GLU pruning is not included since at most 33% sparsity can be achieved.

plexity for language modeling on WikiText-2 test set with sequence length 2048, and on 5-shot accuracy for downstream tasks on the MMLU benchmark (Hendrycks et al., 2020) using the LM Evaluation Harness (Gao et al., 2024).

5.01

6.79

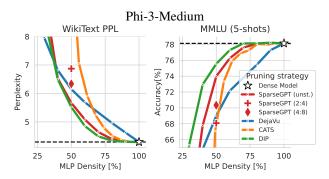


Figure 8. Pareto-curves for WikiText-2 perplexity (left) and 5-shot MMLU accuracy (right) vs MLP density on Phi-3-Medium. The density estimation ignores overhead in DejaVu and SparseGPT.

6.2 Results

DIP+LoRA

Memory footprint We firstly compare DIP against recent static and dynamic pruning methods. In Figure 8 we show the trade-off between model performance and MLP density obtained at varying sparsity levels. Note that certain methods require additional parameters to achieve the sparsification, but we do not account for this in the visualizations, as the resulting overhead depends on the bit-width choice for quantization. DejaVu (Liu et al., 2023) requires sparsity predictors which amount up to 15% of the number of parameters for the dense MLP. For SparseGPT (Frantar & Alistarh, 2023) unstructured and structured pruning, at least 1 extra bit of information for each weight is needed to indicate whether the weight is pruned (Kuzmin et al., 2024).

The memory overhead in this case would be 6.25% and 25% for 16-bit and 4-bit quantized models, respectively. We observe how DIP outperforms both static and dynamic pruning methods at all sparsity levels, even excluding this overhead. Noticeably, we can run Phi-3-Medium in 70% of its full memory footprint, at the cost of only 0.20 perplexity and 0.01% accuracy with respect to the dense model. In comparison, CATS (Lee et al., 2024) incurs in a 0.50 increase in perplexity and 0.88% decrease in accuracy, while for DejaVu (Liu et al., 2023) that is 1.03 higher perplexity and 4.34% lower accuracy. We include in Appendix C additional results for Phi-3-Mini, Llama-v3-8B, and Mistral-7B.

66.57

60.94

59.58

In Table 1 we report more extensive numerical results on four modern LLMs. We consider 50% MLP sparsity for each method, and report the perplexity on WikiText-2 as well as the 5-shot accuracy on the MMLU benchmark. Results for different operating points can be found in Table 3 and Table 4 provided in Appendix C.

Static pruning methods (Frantar & Alistarh, 2023), particularly in the semi-structured setup, show a clear gap in accuracy with respect to the dense model. Gate Pruning and Up Pruning are simple baselines introduced in Section 3. Pruning activations based on the output of the *gate* activations $\sigma(\mathbf{W}^g\mathbf{x})$ (Gate Pruning) severely affects performance on LLama-v3-8B and Mistral-7B. Noticeably, pruning gate layers based on the output of the up component $\mathbf{W}^u\mathbf{x}$ (Up Pruning) instead gives better performance across all tested architectures. Still, both approaches are far from the ideal results achievable with an oracle for GLU Pruning. This oracle shows that, with perfect predictions, model performance would be marginally worse than the dense baseline.

CATS (Lee et al., 2024) is an instance of per-layer threshold pruning on the gate activations. Comparably with Gate Pruning, the method struggles with the Llama-v3-8B and

Table 2. Comparison of throughput for dynamic sparsity methods. We report the highest throughput achieved at a 0.2 and 0.5 increase in perplexity on WikiText-2 over the dense model. The available DRAM capacity is set to fit approximately 50% of each dense model. All models are 4-bit quantized.

	Phi3Med	Phi3Mini	Llama8B	Mistral7B				
Model size	7.4 GB	2.4 GB	4.3 GB	3.9 GB				
DRAM size	4.0 GB	1.5 GB	2.5 GB	2.0 GB				
	Throughput [tok/s] (†)							
Dense	0.29	1.15	0.57	0.54				
	Throughput [tok/s] @ + 0.2 PPL (↑)							
GLU Pruning	0.41	1.59	0.78	0.80				
Up Pruning	0.43	1.69	$\overline{0.76}$	1.03				
CATS	0.40	1.54	0.73	0.69				
DIP	0.43	1.69	0.73	0.98				
DIP-CA	<u>0.45</u>	<u>1.77</u>	<u>0.78</u>	<u>1.08</u>				
	Throughput [tok/s] @ + 0.5 PPL (†)							
GLU Pruning	0.45	1.63	0.84	0.87				
Up Pruning	0.52	1.98	0.87	1.32				
CATS	0.47	1.84	0.83	0.77				
DIP	0.50	1.98	0.82	1.23				
DIP-CA	<u>0.56</u>	<u>2.09</u>	<u>0.89</u>	<u>1.39</u>				

Mistral-7B models. We expect that modifying the CATS approach to use up layer activation might improve its effectiveness. Introducing LoRA adapters over CATS is necessary to boost performance. In contrast, our proposed DIP method without adapters outperforms CATS+LoRA as well as previous baselines on all models and evaluation scenarios. DIP+LoRA further improves overall results, with more pronounced benefits on perplexity scores.

Model throughput In Table 2 we investigate what throughput can be achieved with different sparsity methods and models. For this analysis, all models are quantized to INT4 using GPTQ (Frantar et al., 2022), and all methods rely on an LFU cache allocated as explained in Section 6.1. We consider a use-case where approximately half of the dense model can fit in DRAM and include the selected DRAM sizes compared against the model sizes in the first part of the table. To numerically evaluate the trade-off between model accuracy and efficiency, we report the highest throughput achieved by each method under fixed operating points in perplexity loss. We first notice that GLU Pruning underperforms Up Pruning, as GLU Pruning needs to load the full gate and up layers to compute GLU activations, which do not fit entirely in DRAM. CATS (Lee et al., 2024) falls short from the previous methods, possibly because of the sub-optimality of Gate Pruning described in Table 1. By itself, DIP reaches similar or marginally worse throughput than Up Pruning. Finally, we assess the benefits of applying Cache-Aware masking to DIP. For all models and perplexity

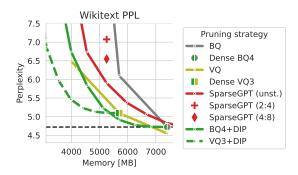


Figure 9. Comparison against quantization and pruning, integration of DIP with 4 bit BQ and 3 bit VQ models on Phi-3-Medium.

operating points, DIP-CA outperforms previous methods. Noticeably, DIP-CA enables running Phi-3-Medium with 46% lower DRAM footprint and 40%, 55% and 93% faster throughput at a perplexity loss of 0.1, 0.2 and 0.5, respectively. Similar improvements are observed for Phi-3-Mini, LLama-v3-8B and Mistral-7B.

6.3 Comparison and compatibility with quantization and static pruning

Quantization is an alternative approach to improve model efficiency. We compare DIP with Blockwise Quantization (BO) (Frantar et al., 2022) and Vector Quantization (VO) (van Baalen et al., 2024) at 2, 3 and 4 bits, as well as static pruning with SparseGPT (Frantar & Alistarh, 2023). Note that in these results we do account for 1 bit of overhead per weight introduced by the pruning mask. In Figure 9 we report results on the comparison and combination of DIP with quantization. DIP applied to 4 bit BQ models compares favorably against both pruning and quantization approaches, reaching 15% lower perplexity than SparseGPT given 4GB of available DRAM. At higher quantization levels, BQ4+DIP outperforms more aggressive BQ quantization, which suggests that increased dynamic sparsity is preferred over decreasing quantization bit-width. VQ3+DIP improves over BQ4+DIP at higher sparsity, similarly outperforming 2-bit VQ, and showing that DIP can be considered as complementary to quantization methods to enhance model efficiency.

6.4 Cache eviction policies

Cache-aware masking By analyzing the distribution of GLU activations in Phi-3-Medium (see Figure 10, left), we observe that a small percentage of activations are orders of magnitude larger than the rest, and will have a bigger impact in determining the layer output. On the other hand, activations between the 30^{th} and 80^{th} percentiles are within one order of magnitude, which means they will have a relatively similar contribution to the linear layer output. DIP-CA

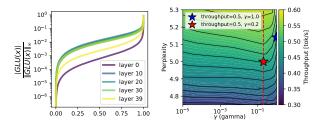


Figure 10. (Left) Distribution of normalized GLU activations for Phi-3-Medium layers. (Right) Effects of changing the DIP-CA γ hyper-parameter on throughput at different levels of perplexity.

exploits this behavior by re-weighting the scores for activations depending on whether they are currently available in cache. This effectively re-orders the activation scores before the TopK selection. In Figure 10 (right), we show the effect of changing the γ hyper-parameter in DIP-CA. For all perplexity levels, γ ranges between 0.1 and 0.3 result in optimal throughput. We include an example by comparing two models achieving the same throughput. The DIP instance with $\gamma=1$ (in blue) has a perplexity of 5.14, while the DIP-CA counterpart with $\gamma=0.2$ reaches lower perplexity of 5.0. Since the activation distributions are comparable across layers (with first and last layers having slightly longer tails), we use a fixed value of γ for all layers.

Comparison with cache eviction policy We compare in Figure 11 different caching policies to handle allocation and eviction of weights in DRAM. Without DRAM caching, inference inference throughput is severly degraded, as the MLP weights need to be loaded from Flash for each token. LFU and LRU policies regulate the eviction of parameters once the DRAM cache is full, by considering how often or recently a neuron in the cache has been used. Results are very similar, with LFU being marginally better. We then observe how even the oracle Belady cache eviction policy does not yield much improvement over practical policies. The throughput for all caching eviction policies overlap under 5.25 perplexity, corresponding to the sparsity level under which the parameter size for the active neurons becomes larger than the DRAM size allocated for caching. In this case, the cache will always contain current neurons, or part of them. In contrast, DIP-CA with a simple LFU policies largely outperforms DIP with all cache eviction policies, even including the optimal Belady algorithm. The improvement also holds at lower sparsity regimes, where part of the active neurons are directly loaded to the processing unit, as there is no remaining space to allocate them in cache.

7 CONCLUSION

This paper addresses the challenge of efficient LLM inference on DRAM constrained mobile devices. Existing

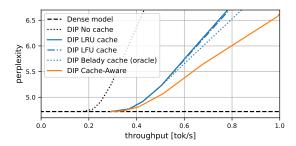


Figure 11. Comparison between cache eviction policies and cacheaware masking. LFU with DIP-CA outperforms all policies including Belady's oracle throughput at a given perplexity threshold.

dynamic sparsity approaches are not effective for modern LLMs which use SwiGLU instead of ReLU, and hence do not exhibit natural sparsity. Instead, we introduce a novel dynamic input pruning (DIP) solution which removes the need for predictors and requires only minimal fine-tuning. DIP can be further aided by a lightweight LoRA adapter to recover part of the accuracy lost during sparsification. Finally, we introduced a cache-ware variant (DIP-CA), which significantly enhances model throughput and reduces memory footprint. Specifically, we achieved a 46% reduction in memory usage and a 40% increase in throughput with respect to dense Phi-3-Medium, with a minimal impact on perplexity. Underscoring these results, we show extensive comparisons against pruning and quantization baselines, as well as ablation studies considering different hardware specifications.

Future work will focus on further refining our sparsity techniques and exploring their applicability to a broader range of models and hardware configurations, including on-device simulation of model performance.

REFERENCES

Abdin, M., Jacobs, S. A., Awan, A. A., Aneja, J., Awadallah, A., Awadalla, H., Bach, N., Bahree, A., Bakhtiari, A., Behl, H., et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.

Akhauri, Y., AbouElhamayed, A. F., Dotzel, J., Zhang, Z., Rush, A. M., Huda, S., and Abdelfattah, M. S. Shadowllm: Predictor-based contextual sparsity for large language models. *arXiv preprint arXiv:2406.16635*, 2024.

Alizadeh, K., Mirzadeh, I., Belenko, D., Khatamifard, K., Cho, M., Del Mundo, C. C., Rastegari, M., and Farajtabar, M. Llm in a flash: Efficient large language model inference with limited memory. *arXiv preprint arXiv:2312.11514*, 2023.

Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., Fan,

- Y., Ge, W., Han, Y., Huang, F., et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- Belady, L. A. A study of replacement algorithms for a virtual-storage computer. *IBM Systems journal*, 5(2): 78–101, 1966.
- Counterpoint. Smartphones beat dram drum to meet performance demand, 2021. URL https://www.counterpointresearch.com/insights/smartphones-dram-trends-2019-2020.
- Danielson, S. How much ram should a phone have, 2023. URL https://robots.net/tech/how-much-ram-should-a-phone-have.
- Dauphin, Y. N., Fan, A., Auli, M., and Grangier, D. Language modeling with gated convolutional networks. In *International conference on machine learning*, pp. 933–941. PMLR, 2017.
- Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al. The llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024.
- Elfwing, S., Uchibe, E., and Doya, K. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. *Neural networks*, 107:3–11, 2018.
- Frantar, E. and Alistarh, D. Sparsegpt: Massive language models can be accurately pruned in one-shot. In *International Conference on Machine Learning*, pp. 10323–10337. PMLR, 2023.
- Frantar, E., Ashkboos, S., Hoefler, T., and Alistarh, D. Gptq: Accurate post-training quantization for generative pretrained transformers. *arXiv preprint arXiv:2210.17323*, 2022.
- Gao, L., Tow, J., Abbasi, B., Biderman, S., Black, S., DiPofi, A., Foster, C., Golding, L., Hsu, J., Le Noac'h, A., Li, H., McDonell, K., Muennighoff, N., Ociepa, C., Phang, J., Reynolds, L., Schoelkopf, H., Skowron, A., Sutawika, L., Tang, E., Thite, A., Wang, B., Wang, K., and Zou, A. A framework for few-shot language model evaluation, 07 2024. URL https://zenodo.org/records/12608602.
- Gholami, A., Yao, Z., Kim, S., Hooper, C., Mahoney, M. W., and Keutzer, K. Ai and memory wall. *IEEE Micro*, 2024.
- Grossi, T. and Lorenz, J. Smartphone memory: Gen ai upgrades to drive spike in dram demand, 2024. URL https://www.yolegroup.com/technology-outlook/smartphone-memory-gen-ai-upgrades-to-drive-spike-in-dram-demand.

- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding. arXiv preprint arXiv:2009.03300, 2020.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Huang, W., Liu, Y., Qin, H., Li, Y., Zhang, S., Liu, X., Magno, M., and Qi, X. Billm: Pushing the limit of post-training quantization for llms. *arXiv preprint arXiv:2402.04291*, 2024.
- Jain, A. and Lin, C. Back to the future: Leveraging belady's algorithm for improved cache replacement. *ACM SIGARCH Computer Architecture News*, 44(3):78–89, 2016.
- Jhang, C.-J., Xue, C.-X., Hung, J.-M., Chang, F.-C., and Chang, M.-F. Challenges and trends of sram-based computing-in-memory for ai edge devices. *IEEE Transactions on Circuits and Systems I: Regular Papers*, 68: 1773–1786, 2021.
- Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., Casas, D. d. l., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., et al. Mistral 7b. arXiv preprint arXiv:2310.06825, 2023.
- Kuzmin, A., Nagel, M., Van Baalen, M., Behboodi, A., and Blankevoort, T. Pruning vs quantization: Which is better? *Advances in neural information processing systems*, 36: 62414–62427, 2023.
- Kuzmin, A., Nagel, M., Van Baalen, M., Behboodi, A., and Blankevoort, T. Pruning vs quantization: which is better? *Advances in neural information processing systems*, 36, 2024.
- Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C. H., Gonzalez, J., Zhang, H., and Stoica, I. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pp. 611–626, 2023.
- Labaran, A. Apple nyme vs ufs 4.0 storage, 2023. URL https://www.allroundreview.com/practical-differences-between-the-apple-nyme-vs-ufs-40-storage.
- Lee, J.-Y., Lee, D., Zhang, G., Tiwari, M., and Mirhoseini, A. Cats: Contextually-aware thresholding for sparsity in large language models. *arXiv preprint arXiv:2404.08763*, 2024.

- Liu, W., Cui, J., Li, T., Liu, J., and Yang, L. T. A space-efficient fair cache scheme based on machine learning for nyme ssds. *IEEE Transactions on Parallel and Distributed Systems*, 34(1):383–399, 2022.
- Liu, Z., Wang, J., Dao, T., Zhou, T., Yuan, B., Song, Z., Shrivastava, A., Zhang, C., Tian, Y., Re, C., et al. Deja vu: Contextual sparsity for efficient llms at inference time. In *International Conference on Machine Learning*, pp. 22137–22176. PMLR, 2023.
- Lykouris, T. and Vassilvitskii, S. Competitive caching with machine learned advice. *Journal of the ACM (JACM)*, 68 (4):1–25, 2021.
- Ma, X., Fang, G., and Wang, X. Llm-pruner: On the structural pruning of large language models. *Advances in neural information processing systems*, 36:21702–21720, 2023.
- Mirzadeh, I., Alizadeh, K., Mehta, S., Del Mundo, C. C., Tuzel, O., Samei, G., Rastegari, M., and Farajtabar, M. Relu strikes back: Exploiting activation sparsity in large language models. arXiv preprint arXiv:2310.04564, 2023.
- Nagel, M., Fournarakis, M., Amjad, R. A., Bondarenko, Y., Van Baalen, M., and Blankevoort, T. A white paper on neural network quantization. arXiv preprint arXiv:2106.08295, 2021.
- Naveed, H., Khan, A. U., Qiu, S., Saqib, M., Anwar, S., Usman, M., Akhtar, N., Barnes, N., and Mian, A. A comprehensive overview of large language models. *arXiv* preprint arXiv:2307.06435, 2023.
- Roumeliotis, K., Tselikas, N., and Nasiopoulos, D. Llama 2: Early adopters' utilization of meta's new open-source pretrained model. preprints2023 2023072 142. *Eng.*). *DOI*, 10, 2023.
- Sadek, K. A. A. and Elias, M. Algorithms for caching and mts with reduced number of predictions. In *The Twelfth International Conference on Learning Representations*, 2024.
- Shah, I., Jain, A., and Lin, C. Effective mimicry of belady's min policy. In 2022 IEEE International Symposium on High-Performance Computer Architecture (HPCA), pp. 558–572. IEEE, 2022.
- Soboleva, D., Al-Khateeb, F., Myers, R., Steeves, J. R., Hestness, J., and Dey, N. SlimPajama: A 627B token cleaned and deduplicated version of RedPajama. https://cerebras.ai/blog/slimpajama-a-627b-token-cleaned-and-deduplicated-version-of-redpajama, June 2023. URL https://huggingface.co/datasets/cerebras/SlimPajama-627B.

- Song, C., Han, X., Zhang, Z., Hu, S., Shi, X., Li, K., Chen, C., Liu, Z., Li, G., Yang, T., et al. Prosparse: Introducing and enhancing intrinsic activation sparsity within large language models. arXiv preprint arXiv:2402.13516, 2024a.
- Song, Y., Mi, Z., Xie, H., and Chen, H. Powerinfer: Fast large language model serving with a consumer-grade gpu. *arXiv preprint arXiv:2312.12456*, 2023.
- Song, Y., Xie, H., Zhang, Z., Wen, B., Ma, L., Mi, Z., and Chen, H. Turbo sparse: Achieving llm sota performance with minimal activated parameters. *arXiv preprint arXiv:2406.05955*, 2024b.
- Sun, M., Liu, Z., Bair, A., and Kolter, J. Z. A simple and effective pruning approach for large language models. *arXiv* preprint arXiv:2306.11695, 2023.
- van Baalen, M., Kuzmin, A., Nagel, M., Couperus, P., Bastoul, C., Mahurin, E., Blankevoort, T., and Whatmough, P. Gptvq: The blessing of dimensionality for llm quantization. *arXiv preprint arXiv:2402.15319*, 2024.
- van der Ouderaa, T. F., Nagel, M., Van Baalen, M., Asano, Y. M., and Blankevoort, T. The llm surgeon. *arXiv* preprint arXiv:2312.17244, 2023.
- Wikipedia. Apple silicon Wikipedia, the free encyclopedia, 2024a. URL https://en.wikipedia.org/w/index.php?title=Apple_silicon&oldid=1250932871. [Online; accessed 14-October-2024].
- Wikipedia. List of qualcomm snapdragon systems on chips Wikipedia, the free encyclopedia, 2024b. URL https://en.wikipedia.org/w/index.php?title=List_of_Qualcomm_Snapdragon_systems_on_chips&oldid=1250680585. [Online; accessed 15-October-2024].
- Wikipedia. Flash memory Wikipedia, the free encyclopedia, 2024. URL https://en.wikipedia.org/w/index.php?title=Flash_memory&oldid=1251065401. [Online; accessed 31-October-2024].
- Xue, Z., Song, Y., Mi, Z., Chen, L., Xia, Y., and Chen, H. Powerinfer-2: Fast large language model inference on a smartphone. *arXiv preprint arXiv:2406.06282*, 2024.
- Zhang, S. and Papyan, V. Oats: Outlier-aware pruning through sparse and low rank decomposition. *arXiv* preprint arXiv:2409.13652, 2024.
- Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M., Chen, S., Dewan, C., Diab, M., Li, X., Lin, X. V., et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.

Efficient LLM Inference using Dynamic Input Pruning and Cache-Aware Masking

Zhao, W. X., Zhou, K., Li, J., Tang, T., Wang, X., Hou, Y., Min, Y., Zhang, B., Zhang, J., Dong, Z., et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.

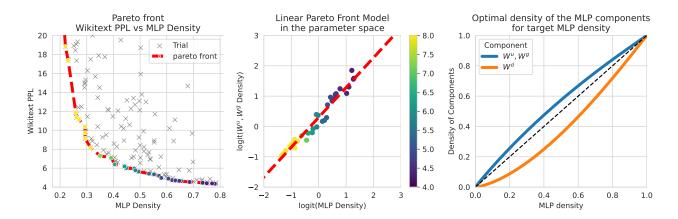


Figure 12. Optimal allocation of weight density between W^d and (W^u, W^g) .

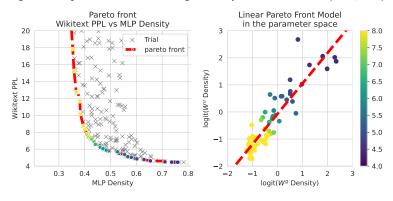


Figure 13. Allocation of weight density between $\widetilde{\boldsymbol{W}}^u$ and $\widetilde{\boldsymbol{W}}^g$ matrices for $\widetilde{\boldsymbol{W}}^d$ with fixed density of 50% for *phi-3-medium model* on the Wikitext dataset.

A HW SIMULATOR DESCRIPTION

The main parameters needed for simulation are therefore the DRAM capacity, and Flash and DRAM read/write speeds. We notice similar specifications in the processors for latest popular smartphone devices: Apple A18 (Wikipedia, 2024a) and Snapdragon® 8s Gen 3 SM8635 platform (Wikipedia, 2024b). While latest UFS versions can reach up to 5.8 GB/s in bandwitdh (Wikipedia, 2024), the effective reading speeds might vary significantly depending on hardware interfaces and whether sequential or random reads are prevalent (Xue et al., 2024). The same applies for NVMe storage systems adopted by Apple devices (Labaran, 2023). Unless otherwise stated, results throughout the paper are based on simulations for Apple A18 with DRAM I/O speed of 60 GB/s, and Flash read speed of 1 GB/s.

A significant part of the available DRAM in smartphones is usually reserved by the OS or other applications (Danielson, 2023). We consider multiple values for available DRAM capacity to investigate how differently sized LLMs perform in memory-constrained scenarios. We include ablations with varying DRAM capacity (Table 5) and Flash read speed

(Table 6) to assess whether the performance is consistent over several device specifications.

Consistently with Xue et al. (2024), for cache simulation, we statically allocate to DRAM all the layers which do not undergo dynamic pruning, such as attention layers, embeddings, but also KV-cache and for certain methods, auxiliary modules like predictors. Since these layers are needed to process each token, pre-loading them in DRAM is an optimal choice (Xue et al., 2024). We then allocate the remaining DRAM capacity uniformly to all the MLP layers to enable the dynamic loading and caching of MLP weights. We did not find significant improvements when exploring non-uniform cache allocation.

B HYPERPARAMETER TUNING

B.1 Density of up, down and gate matrices

We determined the optimal memory allocation for DIP in a three steps procedure illustrated in Figure 13:

1. Run a 2D optimization on Memory vs Perplexity and

	WikiText-2 (Perplexity ↓)			MMLU (5-shot accuracy ↑)				
	Phi3Med	Phi3Mini	Llama8B	Mistral7B	Phi3Med	Phi3Mini	Llama8B	Mistral7B
Dense	4.29	6.01	6.14	5.25	78.14	70.62	65.30	62.68
GLU Pruning (oracle)	4.35	6.04	6.20	5.26	78.28	70.42	65.00	62.35
SparseGPT (unstructured)	5.08	6.63	6.87	5.46	76.26	66.81	60.91	61.23
Gate Pruning	6.36	7.97	10.22	8.67	74.85	62.14	51.62	55.53
Up Pruning	5.68	7.51	8.50	5.74	75.54	64.67	60.29	59.91
DejaVu	5.69	7.74	8.35	5.87	72.15	59.34	55.30	57.58
CATS	5.91	7.43	9.47	8.89	75.91	64.63	55.11	55.65
CATS+LoRA	4.95	6.65	8.61	6.63	75.79	65.64	53.94	56.44
DIP	4.85	6.4	6.66	5.39	77.29	68.88	63.17	61.19
DIP+LoRA	4.62	6.35	6.63	5.38	77.39	$\overline{68.77}$	$\overline{62.95}$	61.50

Table 3. Experimental results for dynamic sparsity methods at 60% MLP density.

Table 4. Experimental results for dynamic sparsity methods at 40% MLP density.

	WikiText-2 (Perplexity ↓)			MMLU (5-shot accuracy ↑)				
	Phi3Med	Phi3Mini	Llama8B	Mistral7B	Phi3Med	Phi3Mini	Llama8B	Mistral7B
Dense GLU Pruning (oracle)	4.29 4.64	6.01 6.24	6.14 6.52	5.25 5.35	78.14 77.74	70.62 69.23	65.30 64.28	62.68 61.67
SparseGPT (unstructured)	6.53	9.51	9.68	6.92	67.79	53.01	48.51	50.15
Gate Pruning	550.27	496.34	>1000	>1000	29.45	26.17	24.51	24.75
Up Pruning	18.57	63.59	68.15	20.44	52.09	29.94	25.32	30.37
DejaVu	6.82	10.55	11.25	6.83	64.77	49.41	41.5	49.38
CATS	196.11	122.93	>1000	>1000	34.8	27.89	25.25	25.12
CATS+LoRA	6.90	9.56	787.70	20.33	60.05	44.67	23.63	26.31
DIP	6.5	8.66	9.01	6.16	72.95	60.33	53.72	55.92
DIP+LoRA	<u>5.64</u>	<u>7.68</u>	<u>8.71</u>	<u>6.03</u>	72.43	<u>61.14</u>	<u>54.35</u>	<u>56.20</u>

determining the pareto optimal configurations

- 2. Modelling the pareto optimal solutions in the parameter space. Here we considered linear models from the target MLP density and the density of the up and gate matrices in logit space.
- 3. Use the fitted model to determine the optimal allocation for a target MLP density.

Using the same procedure, we determined there is no significant gain in using different sparsity levels for the Up W^u and Gate W^g matrices (Figure 13). We found that the optimal allocation is consistent across the tested LLMs.

C ADDITIONAL RESULTS

We show in Figure 14 additional results in terms of perplexity and MLP sparsity for Phi-3-Mini, Llama-v3-8B and Mistral-7B, complementing the study on Phi-3-Medium presented in Section 6.2. Table 4 and Table 3 report the perplexity and accuracy values evaluated at 60% and 40% MLP density respectively. Consistently to the results in Table 1, we don't consider memory overhead introduced by different

methods. The MLP density for CATS and CATS+LoRA may vary up to 2% from the operating point since the fixed estimated threshold may result in slightly different sparsity levels when evaluated on different datasets.

D ABLATION ON HARDWARE SPECIFICATIONS

Table 5. Comparison of throughput for dynamic sparsity methods at different DRAM sizes with Phi-3-Medium quantized to 4 bits. We report the highest throughput achieved at a 0.5 increase in perplexity on WikiText-2 over the dense model.

DRAM size	2 GB	4 GB	6 GB
	Thi	roughput [tol	x/s] ↑
Dense	0.19	0.29	0.71
GLU Pruning	0.24	0.45	1.83
Up Pruning	0.27	0.52	1.35
CATS	0.25	0.47	1.21
DIP-CA	0.31	<u>0.56</u>	<u>1.94</u>

For main experiments we considered a fixed hardware set-

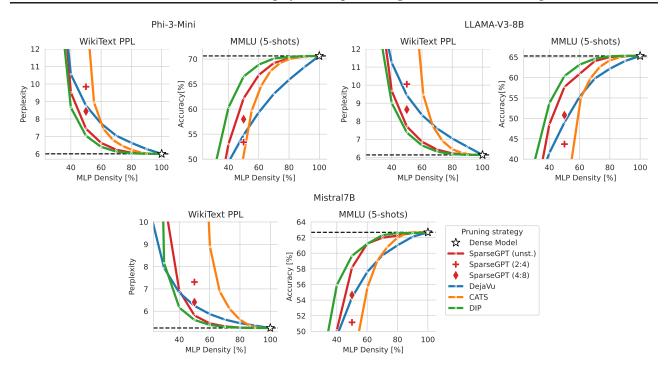


Figure 14. Downstream WikiText Perplexity and MMLU accuracy as a function of MLP Density for Phi-3-Mini, LLama-v3-8B, and Mistral-7B models.

ting in line with the specifications of Apple A18 processors. We now consider different target devices and scenarios by simulating changes in DRAM size and Flash read speed.

DRAM size In Table 5 we consider use-cases with lower DRAM availability at 2GB, simulating the use-case of budget smartphones, and higher DRAM availability at 6GB, which represents high-end devices, or scenarios where the OS and background applications require less memory. In all scenarios, DIP yields better throughput at a fixed perplexity increase of 0.5 over the baseline. Remarkably, DIP improves throughput by 170% against the dense model with 6GB of DRAM. This stems from the increased availability in caching space for the linear layer, which increases the cache hit rate to 89%, compared to the 53% at 4GB and only 8% at 2GB. At even higher DRAM sizes, where the dense model fully fits in cache, we expect GLU Pruning to outperform DIP-CA, as the overhead in loading more weights from DRAM (for GLU Pruning) has a minor impact with respect to the loss in accuracy resulting from sparsifying all MLP layers (in DIP-CA).

Flash reading speed We evaluate DIP against previous methods at different operating points in terms of Flash reading speed. The actual transfer speed can change significantly depending on storage type, hardware interfaces, and data representation for the saved models. While takeaways and relative improvements at different reading speeds do not

Table 6. Comparison of throughput for dynamic sparsity methods at different Flash reading speeds with Phi-3-Medium quantized to 4 bits. We report the highest throughput achieved at a 0.5 increase in perplexity on WikiText-2 over the dense model.

Flash read speed	0.5 GB/s	1 GB/s	2 GB/s
	Thro	ughput [tok/	s] ↑
Dense	0.15	0.29	0.59
GLU Pruning	0.23	0.45	0.91
Up Pruning	0.26	0.52	1.01
CATS	0.24	0.47	0.91
DIP-CA	0.28	<u>0.56</u>	1.09

change, the results in Table 6 highlight how the absolute throughput values change almost at the same rate as the increases in Flash reading speed, confirming that this is the main bottleneck to achieve high-latency in memory-constrained scenarios.