# FlexQuant: Elastic Quantization Framework for Locally Hosted LLM on Edge Devices

Yuji Chai, Mujin Kwen, David Brooks, Gu-Yeon Wei

*Harvard John A. Paulson School Of Engineering And Applied Sciences*

Cambridge, Massachusetts, USA

yuc927@g.harvard.edu

*Abstract*—Deploying LLMs on edge devices presents serious technical challenges. Memory elasticity is crucial for edge devices with unified memory, where memory is shared and fluctuates dynamically. Existing solutions suffer from either poor transition granularity or high storage costs. We propose FlexQuant, a novel elasticity framework that generates an ensemble of quantized models, providing an elastic hosting solution with 15x granularity improvement and 10x storage reduction compared to SoTA methods. FlexQuant works with most quantization methods and creates a family of trade-off options under various storage limits through our pruning method. It brings great performance and flexibility to the edge deployment of LLMs.

*Index Terms*—quantization, edge-devices, elastic serving

## I. INTRODUCTION

Since the introduction of ChatGPT, countless companies in the technology sector have been releasing personal AI assistant products. While OpenAI hosts ChatGPT as a purely in the cloud, these companies are transitioning to a future of AI services that are hosted locally on edge devices, aiming to deliver a truly personalized AI experience that lives entirely on personal devices. Microsoft markets their hybrid hosting concept called "PC + Copilot," in which a language model runs locally on laptops [1]. Apple has integrated Apple intelligence, their personal assistant platform, on virtually their entire new hardware product line [2]. Google has released their Google AI Edge API for Android mobile application developer to deploy Gemini locally on supported Android devices, including their latest smartphone product line [3]–[6]. Several key factors motivate locally hosting LLMs. First, (*Privacy Concerns*) streaming personal information to the cloud for Retrieval Augmented Generation (RAG) greatly increases the risk of privacy breaches [7]. Second, (*Connectivity Issues*) cloud-based inference services rely on a stable internet connection, a privilege that users don't always have on edge devices. Finally, (*Compute Scalability*) the growing compute demand for AI services requires exponential growth in cloud compute capability. Equipping local inference capability could maximize data privacy protection, provide AI service without stable Internet connections, and offload fast-growing compute demand.

Despite the emerging need for local deployment of LLMs on edge devices, it still remains technically challenging. This technical challenge comes from various aspects, including reducing first token latency, improving token generation speed, limiting energy consumption, etc. In this work, we focus our efforts on improving memory footprint elasticity of locally
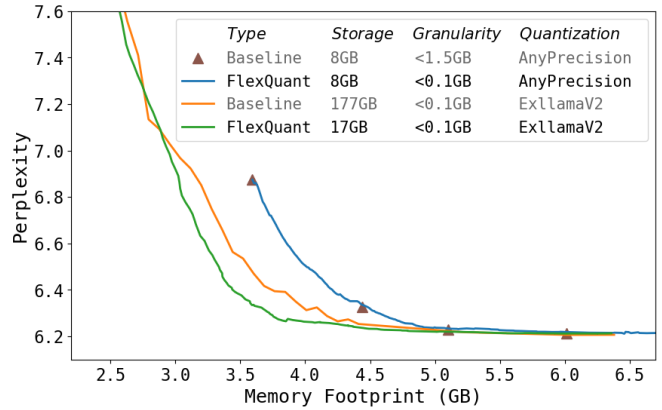


Fig. 1. Perplexity comparison of different elastic hosting methods of quantized Llama 2 7B model. It compares FlexQuant with the baseline method when they are using two quantization methods, ExLlamaV2 and AnyPrecision.

hosted LLMs. These LLMs power Gen-AI services that can be triggered at any time upon a user's request. Because the majority of edge SoCs utilize unified memory architectures [8]–[10], the available memory space is dynamic at the time of a request, depending on memory usage from other applications. Thus, hosting a LLM with a fixed sized memory footprint is not a viable solution and LLM memory elasticity is a necessity.

To enable such elasticity, researchers have designed many compression methods, such as GPTQ [11], ExllamaV2[1], SmoothQuant [12] etc. Those methods, especially ExllamaV2, offer system engineers a flexible design tool that can be used to generate an ensemble of quantized LLMs with various memory footprints. By switching between different models in that ensemble, it serves as an elastic hosting solution. However, this baseline elastic hosting method requires 100GB+ of storage to store all the models. A typical edge device, such as a cell phone, a tablet, or a laptop, usually only has few hundred GB of storage, so a 100GB+ footprint is very costly. A recent work, AnyPrecision [13], provided a better solution where different models could share their parameters resulting in an elastic policy requiring under 10GB of storage. However, it suffers from poor granularity of available model memory footprints and a worse memory-perplexity trade-off compared to ExLlamaV2, as demonstrated in Figure 1. The largest memory footprint gap between two options is larger than

---

[1]https://github.com/turboderp/exllamav2

1.5GB. Based on our analysis in Section II-C, that gap is far from enough for an efficient elastic hosting method. Based on those observations, we propose our elasticity framework FlexQuant with the following contributions:

1) FlexQuant's elastic hosting solution provides a 15x improvement in transition granularity, with the same storage cost, compared to the SoTA elastic hosting methods, as shown in Figure 1.
2) FlexQuant's search algorithm further improves the memory-accuracy trade-off pareto frontier, while achieving more than 10x storage cost reduction and retaining the same level of granularity, as shown in Figure 1.
3) FlexQuant's pruning strategy maintains LLM performance while further reducing the total storage cost of elastic policy by 40%. This method generates a family of elastic hosting options with different storage cost, which provides more flexibility for system designers.

## II. BACKGROUND

### A. Large Language Model Quantization

Reducing the computational cost and memory consumption of LLMs is a pressing challenge, particularly when deploying them on edge devices. Quantization, a method used to approximate higher-precision data using lower-precision representations, has emerged as a standard solution to mitigating this issue. In this paper, we focus primarily on weight quantization methods that generally fall in one of two categories: quantization aware training (QAT) and post training quantization (PTQ). QAT methods simulate quantization during extensive retraining or finetuning phases to mitigate accuracy loss [14], [15]. Although these methods yield superior downstream accuracy, they are both time-consuming and incredibly taxing on computation and memory resources. As this burden is exacerbated with larger models, comparatively lightweight PTQ methods that require little to no training data are generally favored for quantizing LLMs [13], [14], [16]–[18].

One such example is GPTQ, an exemplary one-shot, per-layer quantization method that iteratively quantizes parameters in blocks and uses approximate second-order information to update remaining, yet-to-be-quantized weights, mitigating quantization loss [11]. Other effective PTQ methods include SpQR [19], Activation-Aware Weight Quantization (AWQ) [17], SmoothQuant [12], ZeroQuant [18], AdaRound [20], and SqueezeLLM [16]. In these methods, regardless of weight precision, actual computations usually rely on standard precision formats like as FP16 or FP8. Our elastic framework for hosting quantized LLMs on edge devices is incredibly versatile and can support any PTQ method to mix quantization levels as long as all computations use the same precision. In this paper, we mainly employ ExLlamaV2 for weight quantization. ExLlamaV2 is based on GPTQ and supports mixing quantization levels both across layers and within individual modules, making it particularly well suited for FlexQuant. ExLlamaV2 minimizes quantization error over a diverse calibration set while meeting a target average bitrate for the model. We find

that this PTQ method offers the best trade-off between memory footprint, accuracy, and latency for our experiments.

### B. Elastic Serving

Although PTQ methods are incredibly performant, it is not sufficient to only serve one static model when deploying models on edge devices where many different applications interact with the system, each with its own latency, memory footprint, and accuracy requirements. It would be impossible to meet these varying Service Level Objectives (SLOs) with a single model [21], [22]. However, storing a large portfolio of models of varying sizes is impractical as well. Practical deployment on edge devices demands an elastic solution that can dynamically serve different models in a cost-efficient way.

One state-of-the-art method for efficient deployment is Any-Precision LLM(AnyPrecision) [13]. Starting with a "seed model" quantized (using SqueezeLLM) to the minimum supported bit-width, Any-precision incrementally upscales the model by appending bits to the parameters of the model. At runtime, bits can be truncated or re-appended efficiently, and total storage costs are dictated solely by the largest precision. While this method is elastic, it does not provide a fine-grained transition in memory space as there are large jumps in model footprint between the supported precisions. Furthermore, the method does not identify how to dynamically set precision to meet different SLOs. FlexQuant's versatility allows us to deploy AnyPrecision within the FlexQuant framework, adding granularity to model transitions and providing an ideal trade-off curve between accuracy and memory-footprint.

### C. Transition Granularity

We need to identify an acceptable upper bound on the size of the gaps in memory footprint between models served by the FlexQuant framework. These gaps should be small enough such that we can scale up or scale down model size as memory is requested or freed by other programs on the edge device. It is equally important to minimize memory IO overhead during swapping between different versions of the elastic model. At any point at runtime, we want to identify the model that fits in the current available memory, yet maximizes performance, leaving little "extra" memory, while minimizing memory transition overhead. If we want to swap out a block of parameters, we need to first load the replacement block of parameters in addition to the current hybrid model. We want the size of the additional block of parameters to be small.

One example scenario is deployment on mobile phones where FlexQuant would ideally change model size as apps are opened and closed, suggesting that the ideal granularity is the size of an application. The authors of "End the Senseless Killing" [23] profile the memory allocated by eight popular apps on Android devices, finding that most use around 100MB or less. We additionally profiled application memory allocation over one month of typical usage on an Android phone. Of the top 30 most memory intensive apps in this time period, less than 10% used over 200MB, and the majority of applications used less than 100MB. Thus, an efficient elasticity hosting

method should provide at least 100MB of memory granularity and minimize memory IO overhead.

## III. FlexQuant Framework

FlexQuant enables elastic hosting by generating an ensemble of Elastic Quantization Models ($EQM$). Models in the ensemble have gradually smaller memory footprint to enable memory elasticity. The maximum difference in footprint between two adjacent models defines the ensemble's transition granularity. The total size of the entire ensemble determines its storage overhead. Existing ensemble design methods struggle to generate high granularity memory transition while maintaining low storage cost. FlexQuant automatically identifies $EQM$s to meet this need while improving accuracy-memory trade-off. Its flexible mechanism also allows adaptation to a wide range of quantization methods.

FlexQuant leverages the interchangeability of quantized parameters in different bit-widths. Because an n-bit quantized LLM layer is a numerical approximation of its FP16 counterpart, representations of the layer in different bit-widths still share the same approximation target. Thus, a quantized LLM does not experience significant loss of accuracy if a module is replaced with a slightly lower bit-width counterpart. Empirical observation also supports this intuition. Based on experiments on a Llama 2 7B model, the perplexity drop of an 8-bit model due to replacing one of its module with a 3-bit version's is always under 0.02. This characteristic suggests that a smooth memory-accuracy trade-off, through gradually replacing quantized LLM's layers with their lower bit-width counterparts, is attainable. FlexQuant leverages this characteristic to design a high-granularity $EQM$ ensemble with low storage cost.

### A. EQM Generation Philosophy

Quantization compresses a FP16 model, $M$, into a quantized model, $QM(n)$, with an effective bit-width of $n$. Given two effective bit-widths $n_{low}, n_{up}$ where $n_{low} < n_{up}$, FlexQuant generates a suite of $EQM$s by gradually replacing $QM(n_{up})$'s quantized layers with quantized layers from $QM(n_{low})$. $QM(n_{low})$ and $QM(n_{up})$ define the lower and upper bound of the $EQM$'s memory footprint, respectively. Every $EQM(n_{low}, n_{up})$, defined by $QM(n_{low})$ and $QM(n_{up})$, only uses parameter from those two models and does not require extra storage cost. It could also incorporate additional $m$ quantized models $QM(n_{mid}^i)$, with $n_{low} < n_{mid}^i < n_{up}$ and $0 \le i < m$, to formulate a new set of $EQM(n_{low}, n_{mid}^0, ..., n_{mid}^{m-1}, n_{up})$, for better memory-accuracy trade-off.

The benefits of such a design philosophy are twofold. Firstly, the largest memory footprint difference between two elements in the $EQM(n_{low}, n_{up})$, is no larger than the size of a single LLM module, ensuring high-granularity. Secondly, the high granularity does not require storing new parameters for every element in the $EQM(n_{low}, n_{up})$. No matter how many module replacement trajectories we generate

from $EQM(n_{low}, n_{up})$, they only utilize parameters from $QM(n_{low})$ and $QM(n_{up})$, which ensures low storage cost.

### B. Efficient Design Space Navigation

Although FlexQuant's $EQM$ generation method brings multifold benefits, it still suffers from a major technical challenge: navigating its massive design space. For an $EQM$ ensemble, formulated by $m + 2$ QMs, every element in $EQM(n_{low}, n_{mid}^0, ..., n_{mid}^{m-1}, n_{up})$ has a candidate pool with the size of $(m+2)^{\#module}$, where $\#module$ indicates number of modules in a language model. As a lower bound example, a 2 model $EQM$ ensemble for Llama-2 7B has a design space with the size of $2^{66} = 7.38 * 10^{19}$. To efficiently navigate a design space of this size, FlexQuant narrows down the available options by leveraging deployment constraints. In a realistic deployment scenario, in addition to transition granularity and storage cost, it is also important to limit the transitional memory cost due to module replacement. With this consideration, FlexQuant only allows for module transition from its current bit-with to its lower bit-width counterpart and forbids backward transition. This constraint minimizes the total memory IO cost during continuous transition.

With this design space constraint, the search for the next $EQM^k(n_{low}, n_{mid}^0, ..., n_{mid}^{m-1}, n_{up})$ simplifies into selecting which module to replace with its lower bit-width version. The upper bound of design space's complexity dramatically reduced to $(m + 1) * \#module$. In the case of 2 model $EQM$ ensemble designed for Llama-2 7B, the complexity reduces to 66 for selecting the next configuration. However, this simplification is not enough. The total design space complexity of generating every configuration in the model ensemble, $EQM(n_{low}, n_{mid}^0, ..., n_{mid}^{m-1}, n_{up})$ could be as large as $\prod_{i=1}^{(m+1)*\#module} i$, which is still too large.

We use an approach that is inspired by Monte Carlo Tree Search [24], and the search for $EQM$ generation could be viewed as a tree traversal process. The $EQM$ search starts from the $QM(n_{up})$, which could be viewed as the root of the tree structure. Stemming from the root, there are in total $(m + 1) * \#module$ options as the next configuration. All options would become their own subtree's root and the pattern goes on, until every stem reaches the end. All stems would end with the configuration reaching the $QM(n_{low})$ as leaf nodes. A complete tree traversal path from the root node to a lead node represents a valid $EQM$ suite.

To efficiently navigate the tree, we follow a similar formulation to the Monte Carlo Tree search. FlexQuant only keeps $\#stem$ of the $EQM$ ensemble during search process. At the start of every iteration, FlexQuant loops through every last $EQM_{last}$ in all ensembles. Following the one-way design policy, every successor $EQM$, generated by replacing a module of those $EQM_{last}$ into a lower bit-width counterpart, is added into a candidate pool. To narrow down the scope, it utilizes a pre-generated sensitivity analysis to reserve only the top $\#branch$ candidate for every $EQM_{last}$. The sensitivity analysis is generated by replacing a single module from the $QM(n_{up})$ into its lower bit counterparts

and calculating Euclidean distance between logit outputs from the hybrid model and $QM(n_{up})$. Every module using a bit-width selected from $n_{low}, n_{mid}^0, ..., n_{mid}^{m-1}$, would be associated with a numerical value. The value indicates the additional quantization error resulting from replacing such module with their $n_{up}$ counterpart. The smaller the error, the higher its rank.

After sensitivity analysis, FlexQuant directly evaluates model performance instead of relying on a pre-built sensitivity analysis. To keep the search time reasonable, we utilize the same calibration set used in GPTQ and ExLlamaV2. The calibration set is a diverse collection of token sequences that generally reflect the active numerical range during actual deployment. It serves as a quick evaluation to estimate model performance without having to run a full-scale benchmark. Its effectiveness has been proven in various previous works. FlexQuant's search method estimates the performance of each candidate $EQM$, before completing the entire tree traversal. Based on the evaluation results, top $\#stem$ $EQM$s are added into the ensemble and the search process moves onto the next iteration. The iterative process will continue, until every $EQM_{last}$ in every ensemble has reached $QM(n_{low})$. The ensemble with the best trade-off curve is selected as the final output. The pseudo code for FlexQuant's tree search algorithm is presented in Algorithm 1.

---

**Algorithm 1:** Tree Search for EQM Generation

$ensemble \leftarrow list([QM(n_{up})])$;
$EQM_{last} \leftarrow getLastEQM(ensemble)$;
**while** $EQM_{last} \neq QM(n_{low})$ **do**
    $EQM_{cand} \leftarrow candidateGenerator(EQM_{last})$;
    $EQM_{cand} \leftarrow analysisFilter(EQM_{cand})$;
    $ppl = list()$;
    **for** $eqm \in EQM_{cand}$ **do**
        $ppl.append(calibrationEval(eqm))$;
    **end**
    $EQM_{cand}.rank(ppl)$;
    $EQM_{new} \leftarrow EQM_{cand}[0 : \#stem - 1]$;
    $ensemble.append(EQM_{new})$;
    $EQM_{last} \leftarrow getLastEQM(ensemble)$;
**end**
**return** $ensemble$;

---

Figure 2 illustrates an example with two stems and three branches. At the start of the second iteration, the last $EQM$s from the ensembles are $EQM_{1,0}$ and $EQM_{1,1}$. All of their potential successor $EQM$s are added to the candidate pool. After going through the sensitivity analysis, only $EQM_{2,0}, EQM_{2,1}, ... EQM_{2,5}$ are selected for calibration set evaluation. Only the top two stems, $EQM_{2,2}$ and $EQM_{2,5}$, are selected into the ensemble and become the new $EQM_{last}$ in the search process.

*C. EQM Pruning Strategy*

As mentioned above, it is possible to generate an $EQM$ ensemble with more than 2 starting $QM$s to improve its
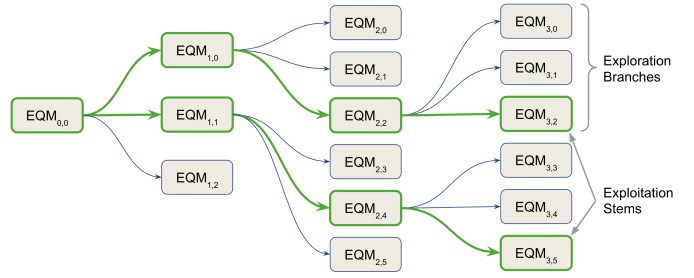


Fig. 2. An example of FlexQuant's tree search for EQM ensemble. The shown search process has two exploitation stems and three exploration branches.

memory and accuracy trade-off. However, not all parameters from the $QM(n_{mid}^i)$ models are necessary nor do they contribute equally to improving the trade-off. Many parameters are not even used in the final $EQM$ ensemble. This observation inspired our pruning strategy that reduces the overhead introduced by $QM(n_{mid}^i)$ models while retaining the accuracy improvements.

FlexQuant's pruning strategy relies on a ranking of the importance or effectiveness of every module introduced by $QM(n_{mid}^i)$. Each module's rank is determined by the number of models in the $EQM$ that leverage its parameters. For modules that have similar amount of usage, ranks are determined by whether the $EQM$ utilizes them at lower footprints. As shown in Figure 1, the accuracy of an $EQM$ is much more sensitive to changes in parameters at smaller memory footprints, shown by the decreasing gradient of the perplexity trade-off curve. Thus, all modules from $QM(n_{mid}^i)$ are ranked first by their overall usage in $EQM$ ensemble and then fine-tuned by their usage for models at lower footprints. Based on this ranking system, FlexQuant can prune the generated $EQM$ ensemble, providing additional storage savings. Additionally, the ranking system allows FlexQuant to prune the ensemble at different rates, granting the user more flexibility in balancing accuracy and storage costs.

## IV. EXPERIMENT

We implement our method using ExLlama on Llama 1 7B, Llama 2 7B, and Llama 3 8B and using AnyPrecision on Llama 1 7B and Llama 2 7B [25]–[27]. For both settings, we generate our $EQM$ ensemble using Exllama's built-in calibration set containing a mixture of data including C4, wikitext, code, and multilingual data. We evaluate this $EQM$ ensemble containing a continuum of hybrid models by measuring both their memory footprint and downstream accuracy. We measure the hybrid model's perplexity on C4 [28], WikiText2 [29], and PTB [30]. In addition to measuring perplexity, we evaluate downstream accuracy on ARC Easy, ARC Challenge [31], Hellaswag [32], PIQA [33], and WinoGrande [34]. These benchmarks cover a diverse range of tasks and domains of applications, which demonstrates our $EQM$s' performances in real deployment scenarios.
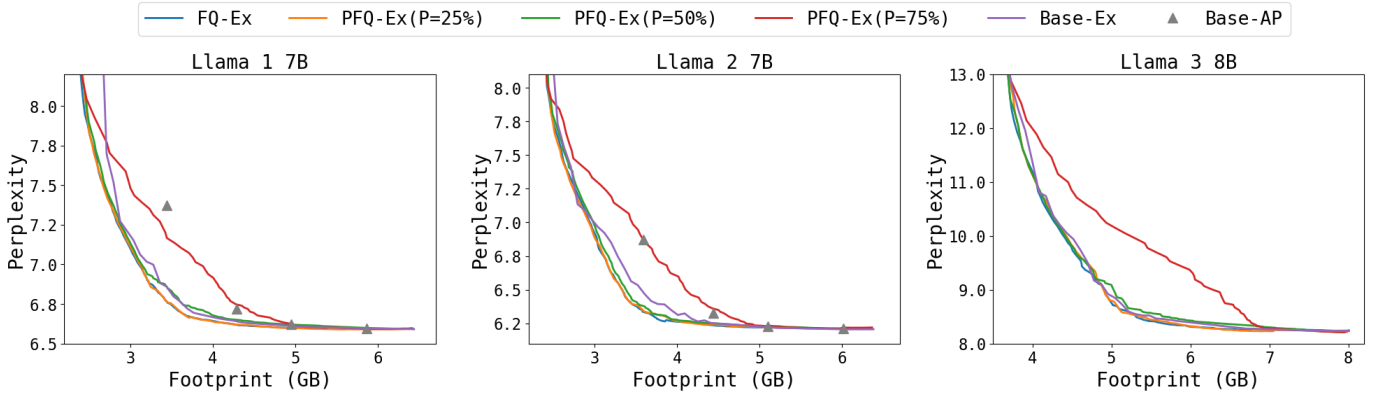
Fig. 3. Perplexity comparison between Base-Ex, Base-AP, FQ-Ex, and PFQ-Ex at different pruning rate. Results for AnyPrecision on Llama 3 8B is omitted due to lack of support in their implementation. The footprint range is different due to differences in parameter count and model architecture.
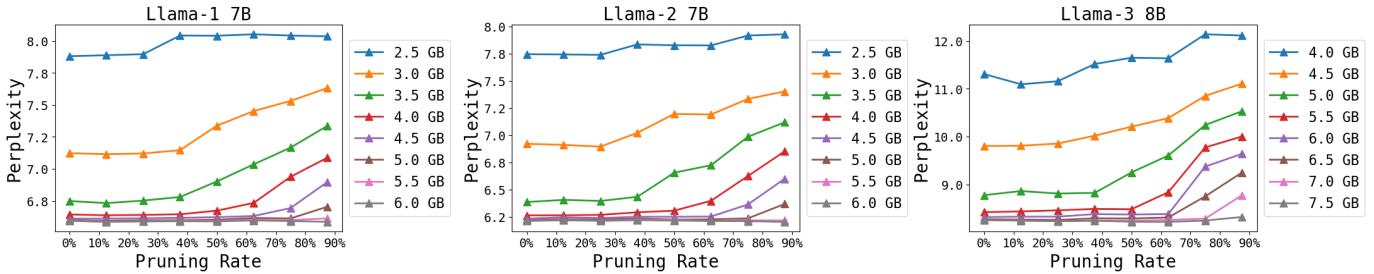


Fig. 4. Perplexity vs pruning rate at varying memory footprint bounds for FQ-Ex on Llama 1, Llama 2 and Llama 3

## V. RESULTS & DISCUSSION

### A. Calibration Set Accuracy Results

Figure 1 shows how FlexQuant, when applied to ExLlamaV2 (FQ-Ex) and to AnyPrecision (FQ-AP), compares to baseline elastic hosting framework using ExLlamaV2 (Base-Ex) and AnyPrecision (Base-AP). We compare these methods by measuring perplexity on the calibration set using Llama 2 7B. Both FQ-Ex and FQ-AP yield smooth transition trajectories with a granularity of ∼100MB. In other words, every consecutive point along each FlexQuant curve represents a potential hybrid model that is only around 100MB larger than the previous one. In comparison, replicating this elastic transition trajectory at 100MB granularity using Base-Ex requires a model ensemble using 177GB in storage, over 10x that of FQ-Ex. Replicating this 100MB granularity with Base-AP is not possible at all. Not only does FQ-Ex require an order of magnitude less storage, it also provides the best memory footprint-accuracy trade-off. Figure 1 shows that FQ-Ex is almost always at the Pareto frontier.

### B. Pruning

As FQ-Ex offers the best trade-off between memory footprint and accuracy, we investigate whether this method can be further improved using pruning. Figure 3 compares Base-Ex, Base-AP, FQ-Ex, and FQ-Ex with pruning (PFQ-Ex). PFQ-Ex is implemented with four pruning rates, P=25%, 50%, and 75%. The four PFQ-Ex curves were all obtained with an $QMs$

of 2.5, 4.0, 5.0, 7.5 bits while varying pruning rates between 0, 0.25, 0.5, ad 0.75, yielding storage costs of 17GB, 15GB, 13GB and 11GB respectively.

Compared to PFQ-Ex, Base-Ex requires over 10x more storage, and Base-AP is roughly 15x less granular. We repeat this experiment using Llama 1 7B, Llama 2 7B, and Llama 3 8B. The memory footprint ranges we explore differ between these three models because Llama 3 8B has more parameters than the other two models. Similarly, the perplexity ranges we achieve differ mainly due to different starting performance for the base models. We find that FQ-Ex, PFQ-Ex with P=25%, and 50% all consistently match or outperform Base-Ex and Base-AP for all three models.

Our pruning method adds more flexibility to FlexQuant. A user can use the same set of quantized models to generate an entire family of trade-off curves with different storage limits by simply adjusting the pruning rate. Starting with the same set of quantized models, FlexQuant can identify multiple candidate hybrid models with similar memory footprints by varying the pruning rate. That is, hybrid models of two different precisions may satisfy the same upper bound depending on how we adjust the pruning rate. Figure 4 demonstrates that models with approximately the same memory footprint can yield different calibration set perplexity depending on the pruning rate at search. The results suggest that relatively low pruning rates maintain similar level of performance as FQ-Ex. However, performance takes a significant hit at around P=40%.
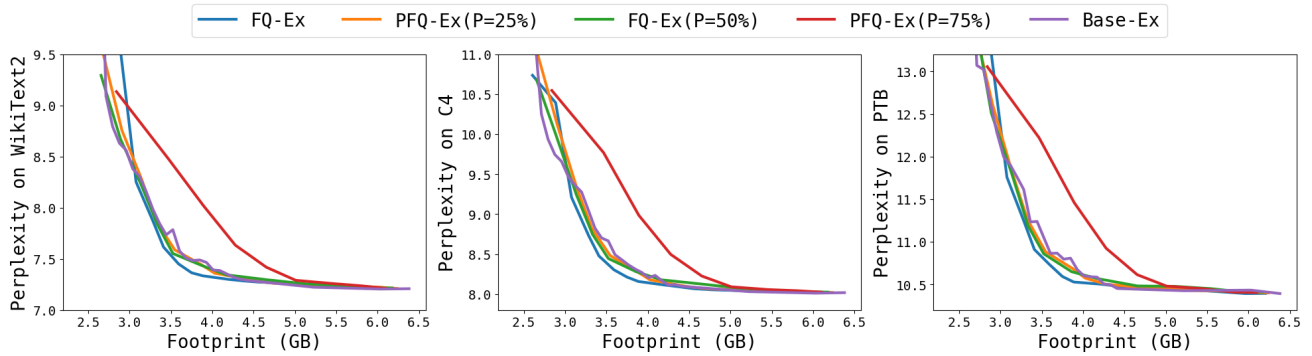
Fig. 5. Downstream perplexity comparison of quantized Llama models between Base-Ex FQ-Ex and PFQ-Ex at different pruning rate.

TABLE I
DOWNSTREAM TASK ACCURACY COMPARISON BETWEEN BASE-EX, FQ-EX AND PFQ-EX WITH 50% PRUNING.

| Task | | $arc_{Challenge}$ | | | $arc_{Easy}$ | | | hellaswag | | | piqa | | | winogrande | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Policy | | Base-Ex | FQ-Ex | PFQ-Ex | Base-Ex | FQ-Ex | PFQ-Ex | Base-Ex | FQ-Ex | PFQ-Ex | Base-Ex | FQ-Ex | PFQ-Ex | Base-Ex | FQ-Ex | PFQ-Ex |
| Storage | | 177 GB | 17 GB | 13 GB | 177 GB | 17 GB | 13 GB | 177 GB | 17 GB | 13 GB | 177 GB | 17 GB | 13 GB | 177 GB | 17 GB | 13 GB |
| Memory Footprint (GB) | 3.0 | 37.40 | -1.40 | -1.40 | 68.98 | +1.73 | +1.73 | 51.47 | -0.60 | -0.80 | 75.63 | +0.15 | +0.54 | 66.61 | +0.24 | -1.10 |
| | 3.5 | 39.42 | +2.04 | +1.08 | 72.05 | +2.57 | +1.77 | 53.21 | +3.02 | +1.87 | 76.22 | +1.14 | +0.60 | 66.69 | +2.22 | +1.34 |
| | 4.0 | 41.38 | +1.19 | +1.11 | 75.01 | +0.35 | +0.20 | 55.81 | +0.76 | +0.55 | 77.04 | +0.76 | +0.71 | 68.27 | +1.10 | +0.40 |
| | 4.5 | 43.17 | -0.26 | -0.43 | 75.46 | +0.00 | -0.25 | 55.81 | +0.93 | +0.77 | 77.75 | +0.26 | +0.00 | 69.37 | +0.32 | -0.39 |
| | 5.0 | 43.17 | +0.17 | -0.25 | 75.46 | +0.13 | -0.17 | 56.85 | +0.08 | +0.03 | 78.01 | +0.11 | +0.12 | 69.38 | +0.39 | -0.16 |
| | 5.5 | 43.34 | +0.17 | -0.26 | 75.58 | +0.51 | +0.01 | 56.92 | +0.20 | +0.07 | 78.12 | +0.07 | +0.06 | 69.38 | +0.37 | +0.23 |
| | 6.0 | 43.34 | +0.34 | +0.08 | 75.58 | +0.51 | +0.01 | 57.03 | +0.14 | -0.01 | 78.18 | +0.05 | +0.11 | 69.45 | +0.44 | +0.16 |

*NOTE: Results of FQ-Ex and PFQ-Ex are presented as accuracy differences, comparing to their Base-Ex counterpart at the same storage footprint.*

## C. Downstream Task Performance

Measuring perplexity on the calibration set alone is not sufficient to confirm model performance. It is important that FlexQuant performs well on downstream tasks that are representative of queries encountered in real-world on-device deployment. We evaluate our method's performance with Llama 2 7B, as shown in Figure 5 and Table I.

Figure 5 compares perplexity on downstream/non-calibration datasets (WikiText2, C4, and PTB) between the ExLlamaV2 baseline and FlexQuant(ExLlamaV2) implemented at four different pruning rates, p=0%, 25%, 50%, and 75%. The top three FlexQuant configurations match or outperform the ExllamaV2 baseline for all tasks. We again see that pruning reduces footprint, but a large pruning rate can destabilize the transition curve and greatly impact performance. Figure 5 demonstrates a dramatic decrease in accuracy when increase pruning rate from 50% to 75%. The observations from Figure 4 and Figure 5 suggest that it is not advisable to apply pruning rates greater than 40–50%.

Table I further compares accuracy on downstream tasks. At 0.5GB intervals, we compare downstream accuracy of the models at that specific memory footprint interval yielded by Base-Ex, FQ-Ex and PFQ-Ex with a 50% pruning rate. For example, the "3.0 GB" row details the downstream task accuracy of the FQ-Ex and PFQ-Ex hybrid models at 3.0 GB compared to the downstream task accuracy of the Base-Ex model directly quantized to 3.0 GB. As shown in the table, the baseline requires 177GB in storage while FQ-Ex and PFQ-Ex only consume 17GB and 13GB of storage respectively. Despite this greater than ten-fold decrease in storage costs, FQ-Ex and PFQ-Ex hybrid models at a given precision typically not only match baseline performance, but also sometimes surpass the downstream task accuracy given by the Base model. This increase in performance is due to our tree search process. The search process is essentially an additional refinement on the quantized LLM, which further optimizes its performance under a footprint limit. As a conclusion, even though our search process is conducted on the calibration dataset, we are presumably not overfitting on this calibration set and sacrificing downstream performance compared to the ExLlamaV2 quantized models.

## VI. CONCLUSION

Growing reliance on LLMs, coupled with the ever-present need for personalization and privacy, has led to increasing interest in locally hosted LLMs. One of the major challenges with locally hosting LLMs is memory elasticity. Existing elastic serving approaches either cannot dynamically adjust model sizes in small increments or fail to scale due to significant storage cost. Our FlexQuant work demonstrates a viable direction to greatly improve transition granularity and storage cost. Future system designers could leverage FlexQuant to have more flexibility and design options when deploying LLMs on edge devices. Additionally, we believe

that a LLM's usage of other key hardware resources, such as energy or compute, also need to be adjusted elastically, a direction we aspire to explore in future works.

REFERENCES

[1] Y. Mehdi, "Introducing Copilot+ PCs - The Official Microsoft Blog," https://blogs.microsoft.com/blog/2024/05/20/introducing-copilot-pcs/.

[2] Apple, "Introducing Apple Foundation Models," https://machinelearning.apple.com/research/introducing-apple-foundation-models, [Accessed 19-11-2024].

[3] "Google AI Edge — Google AI for Developers," https://ai.google.dev/edge, [Accessed 19-11-2024].

[4] K. Semenova, "How to bring your AI Model to Android devices," https://android-developers.googleblog.com/2024/10/bring-your-ai-model-to-android-devices.html, [Accessed 19-11-2024].

[5] T. Darra, "Gemini Nano is now available on Android via experimental access," https://android-developers.googleblog.com/2024/10/gemini-nano-experimental-access-available-on-android.html, [Accessed 19-11-2024].

[6] "Snapdragon X Elite Product Brief," https://docs.qualcomm.com/bundle/publicresource/87-71417-1_REV_F_Snapdragon_X_Elite_Product_Brief.pdf, [Accessed 19-11-2024].

[7] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel *et al.*, "Retrieval-augmented generation for knowledge-intensive nlp tasks," *Advances in Neural Information Processing Systems*, vol. 33, pp. 9459–9474, 2020.

[8] "NVIDIA Jetson AGX Orin — nvidia.com," https://www.nvidia.com/en-us/autonomous-machines/embedded-systems/jetson-orin/, [Accessed 20-11-2024].

[9] A. Brundyn, "Demystifying Unified Memory on Jetson — GTC Digital Spring 2022 — NVIDIA On-Demand — nvidia.com," https://www.nvidia.com/en-us/on-demand/session/gtcspring22-se2600/, [Accessed 20-11-2024].

[10] "Apple introduces M4 Pro and M4 Max," https://www.apple.com/newsroom/2024/10/apple-introduces-m4-pro-and-m4-max/, [Accessed 19-11-2024].

[11] E. Frantar, S. Ashkboos, T. Hoefler, and D. Alistarh, "Gptq: Accurate post-training quantization for generative pre-trained transformers," *arXiv preprint arXiv:2210.17323*, 2022.

[12] G. Xiao, J. Lin, M. Seznec, H. Wu, J. Demouth, and S. Han, "Smoothquant: Accurate and efficient post-training quantization for large language models," in *International Conference on Machine Learning*. PMLR, 2023, pp. 38 087–38 099.

[13] Y. Park, J. Hyun, S. Cho, B. Sim, and J. W. Lee, "Any-precision LLM: Low-cost deployment of multiple, different-sized LLMs," in *Proceedings of the 41st International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, R. Salakhutdinov, Z. Kolter, K. Heller, A. Weller, N. Oliver, J. Scarlett, and F. Berkenkamp, Eds., vol. 235. PMLR, 21–27 Jul 2024, pp. 39 682–39 701. [Online]. Available: https://proceedings.mlr.press/v235/park24e.html

[14] Y. Li, Y. Yu, C. Liang, P. He, N. Karampatziakis, W. Chen, and T. Zhao, "Loftq: Lora-fine-tuning-aware quantization for large language models," *arXiv preprint arXiv:2310.08659*, 2023.

[15] M. Nagel, M. Fournarakis, R. A. Amjad, Y. Bondarenko, M. Van Baalen, and T. Blankevoort, "A white paper on neural network quantization," *arXiv preprint arXiv:2106.08295*, 2021.

[16] S. Kim, C. Hooper, A. Gholami, Z. Dong, X. Li, S. Shen, M. W. Mahoney, and K. Keutzer, "Squeezellm: Dense-and-sparse quantization," *arXiv preprint arXiv:2306.07629*, 2023.

[17] J. Lin, J. Tang, H. Tang, S. Yang, W.-M. Chen, W.-C. Wang, G. Xiao, X. Dang, C. Gan, and S. Han, "Awq: Activation-aware weight quantization for on-device llm compression and acceleration," *Proceedings of Machine Learning and Systems*, vol. 6, pp. 87–100, 2024.

[18] Z. Yao, R. Yazdani Aminabadi, M. Zhang, X. Wu, C. Li, and Y. He, "Zeroquant: Efficient and affordable post-training quantization for large-scale transformers," *Advances in Neural Information Processing Systems*, vol. 35, pp. 27 168–27 183, 2022.

[19] T. Dettmers, R. Svirschevski, V. Egiazarian, D. Kuznedelev, E. Frantar, S. Ashkboos, A. Borzunov, T. Hoefler, and D. Alistarh, "Spqr: A sparse-quantized representation for near-lossless llm weight compression," *arXiv preprint arXiv:2306.03078*, 2023.

[20] M. Nagel, R. A. Amjad, M. Van Baalen, C. Louizos, and T. Blankevoort, "Up or down? adaptive rounding for post-training quantization," in *International Conference on Machine Learning*. PMLR, 2020, pp. 7197–7206.

[21] W. Yin, R. Yi, D. Xu, G. Huang, M. Xu, and X. Liu, "Elms: Elasticized large language models on mobile devices," *arXiv preprint arXiv:2409.09071*, 2024.

[22] W. Yin, M. Xu, Y. Li, and X. Liu, "Llm as a system service on mobile devices," *arXiv preprint arXiv:2403.11805*, 2024.

[23] N. Lebeck, A. Krishnamurthy, H. M. Levy, and I. Zhang, "End the senseless killing: Improving memory management for mobile operating systems," in *2020 USENIX Annual Technical Conference (USENIX ATC 20)*. USENIX Association, Jul. 2020, pp. 873–887. [Online]. Available: https://www.usenix.org/conference/atc20/presentation/lebeck

[24] R. Coulom, "Efficient selectivity and backup operators in monte-carlo tree search," in *Computers and Games*, 2006. [Online]. Available: https://api.semanticscholar.org/CorpusID:16724115

[25] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar *et al.*, "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.

[26] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale *et al.*, "Llama 2: Open foundation and fine-tuned chat models," *arXiv preprint arXiv:2307.09288*, 2023.

[27] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan *et al.*, "The llama 3 herd of models," *arXiv preprint arXiv:2407.21783*, 2024.

[28] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, "Exploring the limits of transfer learning with a unified text-to-text transformer," *Journal of Machine Learning Research*, vol. 21, no. 140, pp. 1–67, 2020. [Online]. Available: http://jmlr.org/papers/v21/20-074.html

[29] S. Merity, C. Xiong, J. Bradbury, and R. Socher, "Pointer sentinel mixture models," *arXiv preprint arXiv:1609.07843*, 2016.

[30] M. P. Marcus, M. A. Marcinkiewicz, and B. Santorini, "Building a large annotated corpus of english: the penn treebank," *Comput. Linguist.*, vol. 19, no. 2, p. 313–330, Jun. 1993.

[31] P. Clark, I. Cowhey, O. Etzioni, T. Khot, A. Sabharwal, C. Schoenick, and O. Tafjord, "Think you have solved question answering? try arc, the ai2 reasoning challenge," *arXiv preprint arXiv:1803.05457*, 2018.

[32] R. Zellers, A. Holtzman, Y. Bisk, A. Farhadi, and Y. Choi, "Hellaswag: Can a machine really finish your sentence?" *arXiv preprint arXiv:1905.07830*, 2019.

[33] Y. Bisk, R. Zellers, J. Gao, Y. Choi *et al.*, "Piqa: Reasoning about physical commonsense in natural language," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 05, 2020, pp. 7432–7439.

[34] K. Sakaguchi, R. L. Bras, C. Bhagavatula, and Y. Choi, "Winogrande: An adversarial winograd schema challenge at scale," *Communications of the ACM*, vol. 64, no. 9, pp. 99–106, 2021.