

SQFT: Low-cost Model Adaptation in Low-precision Sparse Foundation Models

Anonymous ACL submission

Abstract

Large pre-trained models (LPMs), such as large language models, have become ubiquitous and are employed in many applications. These models are often adapted to a desired domain or downstream task through a fine-tuning stage. This paper proposes SQFT, an end-to-end solution for low-precision sparse parameter-efficient fine-tuning of LPMs, allowing for effective model manipulation in resource-constrained environments. Additionally, an innovative strategy enables the merging of sparse weights with low-rank adapters without losing sparsity and accuracy, overcoming the limitations of previous approaches. SQFT also addresses the challenge of having quantized weights and adapters with different numerical precisions, enabling merging in the desired numerical format without sacrificing accuracy. Multiple adaptation scenarios, models, and comprehensive sparsity levels demonstrate the effectiveness of SQFT. **We make SQFT’s fine-tuned models available to reviewers for reproducing our results at: https://anonymous.4open.science/r/sqft_examples-71C7**

1 Introduction

Despite several limitations, such as hallucinations and a significant computational footprint, large pre-trained, foundation, or frontier models have become integral to numerous applications, including language understanding and code generation. These models are trained with extensive corpora on thousands of graphics processing units (GPUs), resulting in outstanding zero-shot performance across various tasks and datasets. However, it is frequently the case that they must be adapted to improve their performance on new tasks or data.

Low-rank adapters (LoRA) (Hu et al., 2022) have demonstrated their effectiveness in model adaptation. However, when LoRA is combined with model compression techniques, e.g., sparsity

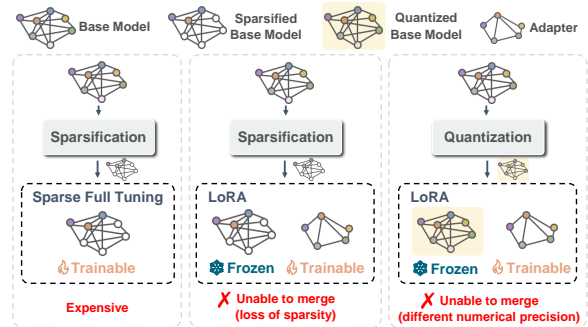


Figure 1: Limitations of existing approaches for fine-tuning sparse and quantized models. Full fine-tuning is expensive. Low-rank adapters (LoRA) for Parameter-efficient Fine-tuning (PEFT) on sparse or quantized models cannot easily merge with the compressed weights due to loss of previously induced sparsity or different numerical precision.

or quantization, several challenges prevent merging these adapters into a single compressed and fine-tuned model, as illustrated in Figure 1. These challenges stem from two primary reasons: **i)** merging dense adapters causes the loss of sparsity in the base model, and **ii)** adapter merging cannot be achieved due to different numerical precisions.

This paper introduces **SQFT**, an end-to-end compression and model adaptation solution for large pre-trained models (LPMs) that alleviates the limitations above. SQFT is designed to sparsify, quantize, and fine-tune large models and can instantiate efficient pipelines that streamline compression techniques. Within the SQFT framework, we propose **Sparse Parameter-Efficient Fine-Tuning (SparsePEFT)**, a strategy to address the adapter merging problem for sparse and quantized model, resulting in more effective high-performing models. Furthermore, SQFT also benefits from weight-sharing techniques applied to traditional parameter-efficient fine-tuning (PEFT) techniques and incorporates insights from state-of-the-art compression techniques. Throughout this paper, we discuss the

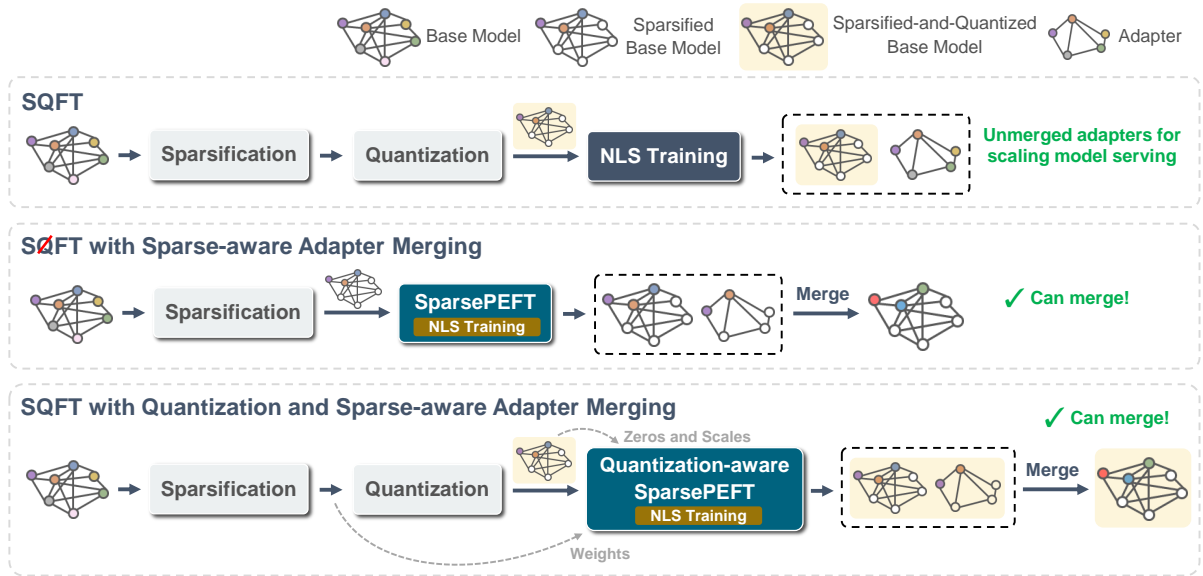


Figure 2: SQFT Overview. Several pipeline configurations can be activated to efficiently fine-tune large models while addressing several limitations of existing approaches.

following contributions:

1. An end-to-end model adaptation solution, SQFT, designed for efficient low-cost configurable pipelines tailored for large pre-trained models with low numerical precision and sparsity.
2. SparsePEFT, a component of SQFT, addresses several limitations in existing parameter-efficient fine-tuning approaches for sparse and quantized models, including the reduction in the cost of fine-tuning, the effective merging of adapters into the sparse model without the loss of sparsity, and the effective merging of components that operate in different numerical precision.
3. Extensive experiments demonstrate the effectiveness of SQFT across different foundation models, sparsity levels and adaptation scenarios.

This paper is organized as follows: Section 2 describes the stages in the proposed end-to-end solution, SQFT. Section 3 discusses SQFT’s evaluation, and we finalize with some concluding remarks in Section 4. Due to page limits, we include a Related Work section, and additional results in the Appendix.

2 Methodology

SQFT fine-tunes large pre-trained models (LPMs) in an efficient multi-stage approach that includes

(1) Sparsification, with an optional reduction in the numerical precision, i.e., Quantization, (2) Fine-tuning with Neural Low-rank Adapter Search (NLS), (3) Sparse Parameter-Efficient Fine-Tuning (SparsePEFT) with optional (4) Quantization-awareness. Figure 2 illustrates the alternative LPM compression and model adaptation pipelines that SQFT can instantiate. In the following sections, we discuss the details of each stage and the benefits of accelerating inference and model serving.

2.1 Sparsification and Quantization Stage

As shown in Figure 2, at the beginning of all possible pipeline configurations, SQFT employs an effective method to induce sparsity in the model. For a given weight matrix $\mathbf{W} \in \mathbb{R}^{m \times n}$, with entries $w_{i,j}$ s.t. $\mathbf{W} = (w_{i,j}), 1 \leq i \leq m, 1 \leq j \leq n$, an arbitrary scoring function, Ψ , is assigned to the proposed solution. This function determines the relative importance of $w_{i,j}$ compared to the other weights in \mathbf{W} . Ψ can be formulated in various ways. For instance, $\Psi(\mathbf{W}) = |\mathbf{W}| \cdot \|\mathbf{X}\|_2$, where \mathbf{X} represents sampled feature input activations, as proposed by Sun et al. (2023). However, it is important to highlight that the proposed end-to-end model fine-tuning solution, SQFT, can utilize any other scoring function. Leveraging the scores from Ψ and a desired level of sparsity, s , we derive the sparsified weight, denoted as \mathbf{W}^p , with a sparsity pattern $S\{\mathbf{W}^p\} = \{(i, j) \mid \mathbf{W}_{i,j}^p \neq 0, 1 \leq i \leq m, 1 \leq j \leq n\}$, s.t. $|S\{\mathbf{W}^p\}| \leq |S\{\mathbf{W}\}|$.

It has been demonstrated that LPMs can tolerate

higher sparsity levels compared with the previous generations of smaller transformer-based models (Frantar and Alistarh, 2023). Our experiments confirm these observations (Section 3). Once SQFT has induced sparsity in the pre-trained weights, W^p enables an optional reduction in their numeric precision. Given the sparsified weights, SQFT applies layer-wise one-shot quantization (Nagel et al., 2020; Frantar et al., 2022a; Wang et al., 2020; Frantar et al., 2022b). Utilizing a selection from state-of-the-art post-training quantization approaches, SQFT identifies the low-precision sparsified weights, denoted as \widehat{W}^p , that given an input X , minimize $\text{argmin}_{\widehat{W}^p} \|W^p X - \widehat{W}^p X\|_2^2$.

Reducing the numerical precision and inducing sparsity on weights frequently decrease the model’s accuracy, requiring fine-tuning to improve performance.

2.2 Fine-tuning with Neural Low-rank Adapter Search (NLS)

Given the sparse quantized weights, \widehat{W}^p , SQFT recovers any drops in accuracy induced by the compression schema and fine-tunes these weights for a specific downstream task. As shown in Figure 2, SQFT employs Neural Low-rank Adapter Search (NLS) (Munoz et al., 2024a) instead of vanilla Low-rank Adapters (LoRA) (Hu et al., 2022), and fine-tunes sparse and quantized model. To justify using NLS, traditional LoRA adapters require assigning the values for several hyperparameters, including their rank r , and the subset of modules where these adapters will be placed. Determining these hyperparameters can be a challenging endeavor. To alleviate this limitation, SQFT extends NLS’ weight-sharing techniques to facilitate the discovery of optimal adapter configurations from a space of elastic adapter configurations. In other words, instead of having a fixed value for the rank, r , we enable elastic configurations, $C = [c_1, \dots, c_n]$, s.t., $r \leftarrow c_i$ depending on the activation of the corresponding sub-adapter.

2.3 SparsePEFT

Fine-tuning the sparse quantized model with adapters effectively improves the model’s performance on a downstream task. However, as illustrated in the middle and right part of Figure 1, a challenge arises when dealing with sparse or quantized weights and dense adapter weights: merging them will i) result in the loss of sparsity on the model’s weights or ii) be unable to merge due

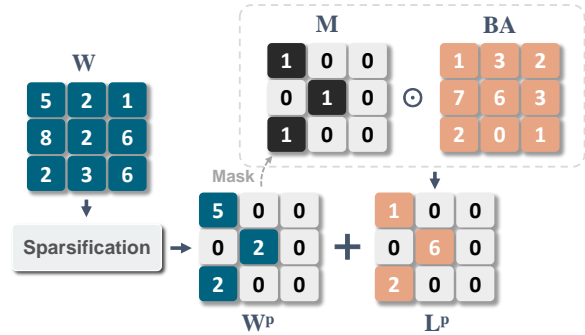


Figure 3: Sparse Parameter-efficient Fine-tuning (SparsePEFT). A binary mask is obtained from the sparsified weights and applied to the adapters, allowing for the later merge without loss of sparsity.

to different numerical precisions. Aiming to address the first limitation, we propose an effective strategy, Sparse Parameter-Efficient Fine-Tuning (SparsePEFT), to make adapters sparsity-aware. As depicted in Figure 3, SparsePEFT applies a binary mask M derived from the initial sparsification of W . This mask is used to sparsify the adapters matrix (denoted as BA) into L^p . The process can be formulated as:

$$L^p = (BA) \odot M, \quad (1)$$

which is activated during the fine-tuning process for sparsity awareness. SparsePEFT enables the merging of the sparsified weights W^p and the adapter weight L^p without sacrificing the sparsity induced early in the compression pipeline as follows,

$$W^p \leftarrow W^p + L^p. \quad (2)$$

In addition to preserving sparsity, SparsePEFT demonstrates comparable (even better) accuracy compared to fine-tuning with dense adapters. Extensive experimental findings substantiate the advantages of SparsePEFT, as detailed in Section 3.

Although SparsePEFT can effectively preserve the model’s sparsity, it presents additional challenges when merging with quantized models, the second limitation we discussed before, which is primarily attributed to the need for the adapter and pre-trained weights to possess identical numerical precision. In the following subsection, we explore a pipeline variation for SQFT that facilitates the integration of sparse quantized weights. This approach aims to address both challenges mentioned above while improving the overall efficiency of the resulting model.

2.4 Quantization-aware SparsePEFT

Building upon the concept of SparsePEFT, we propose Quantization-aware SparsePEFT (QA-SparsePEFT), an extension of SparsePEFT for sparse quantized models. QA-SparsePEFT integrates quantization awareness into SparsePEFT. In most common quantization schemes, the zero point and scales for the target quantized tensor can be determined during the quantization process (e.g., GPTQ (Frantar et al., 2022a)). Within the framework of QA-SparsePEFT, the zeros, and scales of the sparse quantized weights \widehat{W}^p are shared with the adapter. The adapters can be quantized smoothly with the shared fixed zeros and scales, enabling quantization-aware fine-tuning. Formally, given the sparsified pre-trained weight W^p , sparsified adapter weight L^p obtained from SparsePEFT, zeros z and scales s from the quantization of W^p , the quantization process in the proposed QA-SparsePEFT can be formulated as:

$$\widehat{W}_m^p = \text{round} \left(\text{clamp} \left(\frac{(W^p + L^p) - z}{s}, Q_n, Q_p \right) \right), \quad (3)$$

where \widehat{W}_m^p denotes the sparse quantized (merged) weight, $Q_n = -2^{n-1}$ and $Q_p = 2^{n-1} - 1$ (n represents the bit-width of the quantized values). Dequantization is the inverse as follows:

$$\tilde{W}_m^p = \widehat{W}_m^p \times s + z, \quad (4)$$

which applies z and s to approximate W_m^p . Through QA-SparsePEFT, we can obtain the fine-tuned, sparsified low-precision resulting model. Moreover, SQFT with QA-SparsePEFT can run the NLS stage using this schema, which allows us to merge the adapters once an optimal configuration has been discovered.

2.5 Model Serving and Inference Acceleration

Accelerating model serving and inference through sparsification and quantization techniques has shown significant efficacy across various hardware platforms and kernels, demonstrating remarkable speedups. However, for PEFT with a sparsified or quantized model (as shown in Figure 1), the addition of adapter models introduces the computational overhead during inference due to their non-mergeability. SparsePEFT (QA-SparsePEFT) allows adapters to be merged into the sparse (quantized) model, which can reduce adapters' redundancy and computational overhead, leading to more streamlined inference processes. Moreover, quantization techniques further enhance acceleration by

reducing the model size and computational complexity, but balancing the trade-off between acceleration and maintaining competitive accuracy is essential.

In summary, SQFT and its SparsePEFT strategy bring the benefits of adapter merging and maintaining accuracy on sparse or quantization scenarios. The choice between the sparsity level and whether to apply quantization depends on the specific deployment scenario (e.g., task requirements and resource constraints), including the trade-off between model performance, inference speed, and memory efficiency. In the next section, we will delve into further empirical studies to fully understand the strengths and weaknesses of each approach in different settings.

3 Experimental Results

We evaluate SQFT on several state-of-the-art large pre-trained models and datasets. Next, we discuss the setup for our experimental analysis.

3.1 Setup

Models SQFT is evaluated on two state-of-the-art models, including Llama-3-8B¹, Phi-3-Mini-4K-Instruct². To study it more comprehensively, we aim to explore SQFT across different models, scales, and settings.

Datasets and Settings Aligned with other works in the LPMs compression and fine-tuning spaces, SQFT is validated on three experimental settings: **1)** Grade School Math 8K (GSM8K) (Cobbe et al., 2021), **2)** Math reasoning with instruction tuning (following LLM-Adapters (Hu et al., 2023)), including 3 math reasoning datasets: GSM8K, Math Word Problems (MAWPS) (Koncel-Kedziorski et al., 2016), Simple Variations on Arithmetic Math word Problems (SVAMP) (Patel et al., 2021), and **3)** Commonsense reasoning datasets: Boolean Questions (BoolQ) (Clark et al., 2019), Physical Interaction: Question Answering (PIQA) (Bisk et al., 2020), HellaSwag (Zellers et al., 2019), Large-scale Winograd Schema Challenge (WinoGrande) (Sakaguchi et al., 2021), AI2 Reasoning Challenges (Arc-e, Arc-c) (Clark et al., 2018), and Open Book Question Answering (OBQA) (Mihaylov et al., 2018).

The evaluations of our experiments are conducted utilizing *lm-eval-harness* (Gao et al., 2023)

¹<https://huggingface.co/meta-llama/Meta-Llama-3-8B>

²<https://huggingface.co/microsoft/Phi-3-mini-4k-instruct>

Table 1: Results for adapting **Llama-3-8B** to GSM8K. The criterion for mergeable is that there should be no loss in either accuracy or sparsity before and after merging. The evaluation used the default configuration for *lm-eval-harness* (Gao et al., 2023).

Model	Sparsity	Method	Mergeable	Final Precision (Base + Adapter / Base)	GSM8K Test Accuracy(%)	
Llama-3-8B	0%	w/o tune	-	FP16	50.0	
	50%	----- <i>w/o Quantization</i> -----				
		w/o tune	-	FP16	12.5	
		LoRA	✗	FP16 + FP16	50.6	
		Shears	✗	FP16 + FP16	52.2	
		SQFT + SparsePEFT (Ours)	✓	FP16	52.5	
		----- <i>Quantization</i> -----				
		w/o tune	-	INT4	7.0	
		GPTQ + LoRA	✗	INT4 + FP16	48.9	
		SQFT (Ours)	✗	INT4 + FP16	50.0	
SQFT + QA-SparsePEFT (Ours)		✓	INT4	50.2		

in both setting 1 and 3, while following the evaluation from LLM-Adapters in setting 2. We present a comparative analysis of the results obtained from our various pipelines and also compare with vanilla LoRA (Hu et al., 2022), Shears (Munoz et al., 2024a) (a parameter-efficient fine-tuning method for sparse models), and GPTQ + LoRA. For fair comparison, all methods are run in the same environment and with the same configuration. SQFT employs the implementation of Wanda (Sun et al., 2023) as default method for sparsification, and GPTQ in Huggingface³ for quantizing the LPMs and adapters.

Reference configuration Unless stated in the results, we report a reference configuration for SQFT. This configuration is obtained utilizing the heuristic proposed in Munoz et al. (2024b). The heuristic is intuitive and straightforward, activating the configuration with the median of each set of elastic values per module. Spending additional cycles to search the space of configurations might yield even more competitive results, presented in Table 4. Next, we discuss experimental results and studies conducted using SQFT.

3.2 Main Results

3.2.1 Fine-tuning Llama-3 on GSM8K

We begin our evaluation with Llama-3B-8B, assessing its accuracy in a dense mode and after inducing 50% sparsity without fine-tuning on the GSM8K dataset. Subsequently, we execute various pipelines of SQFT. As described in Table 1, for Llama-3-8B at the 50% sparsity level, SQFT recovers the model’s accuracy from 12.5% to 52.5% without employing quantization, while allowing for

³<https://huggingface.co/blog/gptq-integration>

the merging of adapters without sacrificing sparsity (SparsePEFT) and incorporating quantization into the pipeline results in a minor drop in accuracy to 50.2% when enabling the adjustment to merge adapters (QA-SparsePEFT).

More importantly, SQFT with SparsePEFT and QA-SparsePEFT exhibit comparable performance to their corresponding non-mergeable approaches. These results suggest that SQFT with SparsePEFT (QA-SparsePEFT) effectively addresses the limitation of the merging problem encountered when fine-tuning adapters into sparse models (or sparse and quantized models) without any degradation in accuracy. Furthermore, the comparison between LoRA and SQFT with SparsePEFT (or Shears), and between GPTQ + LoRA and SQFT with QA-SparsePEFT without adapter merging, highlights the superior performance of NLS (elastic rank) compared with LoRA (fixed rank). We also explore the performance of a broader range of sparsity levels and conduct more detailed ablation experiments in this experimental setting, which can be found in Sections 3.4 and 3.6, respectively.

3.2.2 Math Reasoning with Instruction Tuning for Phi-3

In addition to Llama-3 on GSM8K, we also investigated the performance of SQFT with the Phi-3 model. Since the Phi-3-series models released by Microsoft are the instruction models currently best-suited for a chat prompt, we evaluate SQFT on three math reasoning datasets for instruction tuning. Table 2 presents the test accuracy for our approaches and baselines. Interestingly, in the full-precision mode (*w/o Quantization*), our proposed SparsePEFT not only achieves the highest average accuracy (77.3%) compared to other approaches

Table 2: Results for **Phi-3-Mini-4K-Instruct** with math instruction tuning. *Mergeable* means that merging the dense adapters with the sparse weights is possible without losing the induced sparsity levels or affecting the desired low numerical precision.

Model	Sparsity Method		Mergeable	Final Precision (Base + Adapter / Base)	Datasets Accuracy(%)			Average
	0%	w/o tune			GSM8K	MAWPS	SVAMP	
Phi-3-Mini-4K-Instruct	0%	w/o tune	-	FP16	64.7	84.5	85.4	78.2
				<i>w/o Quantization</i>				
		w/o tune	-	FP16	38.9	64.7	66.8	56.8
		LoRA	✗	FP16 + FP16	62.5	90.3	77.8	76.9
		Shears	✗	FP16 + FP16	62.3	90.8	76.1	76.4
		SQFT + SparsePEFT (Ours)	✓	FP16	61.9	91.2	78.7	77.3
				<i>Quantization</i>				
		w/o tune	-	INT4	33.4	56.7	64.2	51.4
		GPTQ + LoRA	✗	INT4 + FP16	60.3	89.5	74.8	74.9
		SQFT (Ours)	✗	INT4 + FP16	60.3	90.8	75.6	75.5
	SQFT + QA-SparsePEFT (Ours)	✓	INT4	60.4	90.8	72.9	74.7	

but also uniquely allows for the merging of adapters and sparse weights without any loss of sparsity. This result is achieved without needing an expensive search and by utilizing the heuristic detailed in Section 3.1. However, in quantization mode, the accuracy of SQFT + QA-SparsePEFT (mergeable) is marginally lower compared to the non-mergeable approaches (74.7% vs. 74.9%/75.7%). This result suggests there may be a need to balance the trade-off between accuracy and efficiency. Fortunately, SQFT + QA-SparsePEFT results in a merged fine-tuned quantized model, eliminating the overhead associated with dense adapters.

3.2.3 Fine-tuning Phi-3 on Commonsense Reasoning

Besides the mathematical domain of the first two experimental settings, we also explore SQFT in other areas, e.g., commonsense reasoning. We apply SQFT to fine-tuning the Phi-3 model on a set of unified commonsense training datasets with 83K samples for fine-tuning from BoolQ, PIQA, HellaSwag, WinoGrande, Arc-e, Arc-c, and OBQA. Table 3 compares the test accuracy of the evaluated approaches. SQFT obtains a competitive configuration with Shears, LoRA, and GPTQ + LoRA. However, SQFT has the additional benefit of allowing for the merging without losing the previously induced sparsity, both in full-precision and quantized modes. It is worth noting that SQFT with QA-SparsePEFT shows super competitiveness here, i.e., the most efficient model with high accuracy (among all full-precision and quantized cases).

3.3 Hill-climbing to Better Configurations

The results presented in the previous sections employ the simple heuristic (as detailed in Section

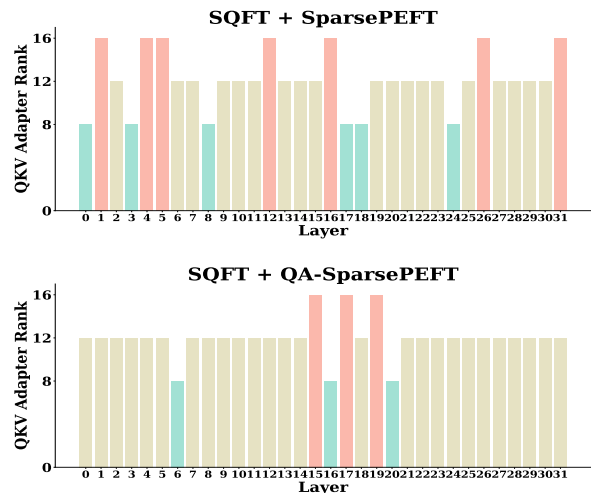


Figure 4: The adapter rank distribution of the optimal configurations obtained from the hill-climbing search algorithm (**Phi-3-Mini-4K-Instruct** with commonsense reasoning).

3.1) to obtain a reference configuration from the NLS search space. However, superior configurations can be discovered with an additional budget. We apply a well-designed hill-climbing search algorithm (Algorithm 1 in Appendix), which starts from the configuration derived from the heuristic and explores its neighboring configurations in a hill-climbing manner based on their validation accuracy. For this purpose, we employed the validation sets from Arc-e, Arc-c, and OBQA, as other datasets do not provide a validation set. As demonstrated in Table 4, a more optimal configuration can be discovered, outperforming the default adapter configuration obtained from the heuristic. Exploring further the search space of elastic adapter ranks produces richer adapter distributions as depicted in Figure 4. More importantly, the test set results re-

Table 3: Results for **Phi-3-Mini-4K-Instruct** with commonsense reasoning. SQFT obtains competitive fine-tuned models with an additional benefit over Shears and LoRA applied to low-precision weights, i.e., SQFT’s adapters can be efficiently merged into the weights without any loss of precision or accuracy. We are reporting a reference submodel for SQFT obtained the heuristic detailed in 3.1, which means that, as shown in Table 4, with an additional cost, SQFT can discover submodels with even higher performance.

Model	SparsityMethod		Mergeable	Final Precision (Base + Adapter / Base)	Datasets Accuracy(%)							Average	
					BoolQ	PIQA	HellaS	WinoG	Arc-e	Arc-c	OBQA		
Phi-3-Mini-4K-Instruct	0%	w/o tune	-	FP16	86.1	80.3	78.5	73.7	83.2	57.5	46.8	72.3	
	w/o Quantization												
		w/o tune	-	FP16	82.5	75.9	69.9	69.1	76.9	50.9	43.4	66.9	
		LoRA	✗	FP16 + FP16	85.6	79.1	75.8	71.5	79.6	53.2	49.4	70.6	
		Shears	✗	FP16 + FP16	85.2	78.9	75.7	72.6	80.1	53.3	50.4	70.9	
		SQFT + SparsePEFT (Ours)	✓	FP16	84.0	78.8	75.5	72.1	80.1	53.5	48.6	70.4	
	Quantization												
		w/o tune	-	INT4	81.4	75.2	68.5	68.2	75.9	50.3	40.2	65.7	
		GPTQ + LoRA	✗	INT4 + FP16	85.3	79.1	75.3	72.5	79.5	54.6	47.2	70.5	
		SQFT (Ours)	✗	INT4 + FP16	85.1	79.0	75.4	71.2	79.6	54.1	48.8	70.5	
		SQFT + QA-SparsePEFT (Ours)	✓	INT4	83.7	80.1	74.1	73.6	80.1	55.1	48.2	70.7	

Table 4: Hill-climbing searching results for **Phi-3-Mini-4K-Instruct** with commonsense reasoning.

Model	Sparsity	Method	Sub-Adapter	Validation Datasets Accuracy(%)				Test Datasets Accuracy(%)							
				Arc-e	Arc-c	OBQA	Average	BoolQ	PIQA	HellaS	WinoG	Arc-e	Arc-c	OBQA	Average
Phi-3-Mini-4K-Instruct	50%	SQFT + SparsePEFT	Heuristic	79.3	50.8	47.4	59.2	84.0	78.8	75.5	72.1	80.1	53.5	48.6	70.4
			Hill-climbing	80.2	51.8	47.6	59.9	84.3	78.9	75.4	72.0	80.1	54.3	49.4	70.6
		SQFT + QA-SparsePEFT	Heuristic	80.0	51.5	45.4	59.0	83.7	80.1	74.1	73.6	80.1	55.1	48.2	70.7
			Hill-climbing	80.4	53.5	46.2	60.0	83.6	79.7	74.1	73.7	80.1	56.2	48.8	70.9

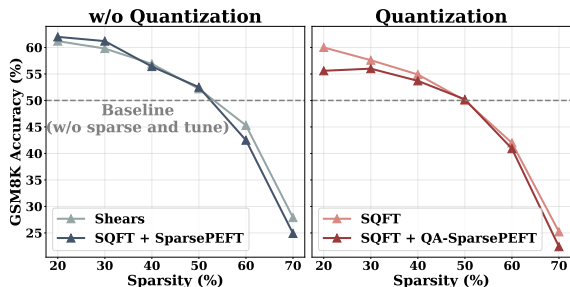


Figure 5: Comparison of various sparsity levels for Llama-3-8B with GSM8K. SQFT achieves similar performance as Shears but with the added benefit of merging adapters with different numerical precision.

veal a significant improvement in the performance of the Arc-c and OBQA datasets, which suggests that an appropriate validation set can assist in identifying the optimal adapter configuration.

3.4 Exploring a Broader Range of Sparsity Levels

All our previous experiments employ 50% sparsity as it is moderate and mild. In this section, we explored the behavior of SQFT in a broader range of sparsity levels. As shown in Figure 5, the model’s accuracy experiences a significant drop between a

sparsity of 60% and 70%. We denote this range as the critical sparsity threshold, representing the boundary at which the model’s performance begins to degrade notably. Through our recovery downstream fine-tuning strategy, models with up to 50% sparsity (even with quantization) can achieve comparable performance with the original dense model (represented by the baseline in the figure) on the downstream task. This 50% sparsity can be defined as the optimal sparsity level, as it represents the point of balance where the model maintains high performance while achieving computational efficiency. Moreover, there is little difference in accuracy between our mergeable approaches and non-mergeable methods, which illustrates the effectiveness of our proposed SparsePEFT.

3.5 Cost Analysis of Pipeline Configurations

The different versions of SQFT’s pipelines incur various costs that allow users to choose based on their fine-tuning budget. Table 6 details the characteristics of each pipeline configuration, e.g., whether we can merge the adapters, the precision of the based model and the adapters, and the cost of each configuration. Two assumptions are

Table 5: Ablation studies for LoRA vs. NLS (**Llama-3-8B** with GSM8K). Compared to LoRA, NLS demonstrates significantly better accuracy performance across all possible pipelines of SQFT and different sparsity levels.

Model	Sparsity	Method	Mergeable	Final Precision (Base + Adapter / Base)	Fine-tune Approach	GSM8K Test Accuracy(%)
Llama-3-8B	30%	Shears	✗	FP16 + FP16	LoRA	58.2
					NLS	59.8^{+1.6}
		SQFT + SparsePEFT (Ours)	✓	FP16	LoRA	60.0
					NLS	61.2^{+1.2}
		SQFT (Ours)	✗	INT4 + FP16	LoRA	56.7
					NLS	57.6^{+0.9}
		SQFT + QA-SparsePEFT (Ours)	✓	INT4	LoRA	54.8
					NLS	56.0^{+1.2}
		Shears	✗	FP16 + FP16	LoRA	50.6
					NLS	52.2^{+1.6}
		SQFT + SparsePEFT (Ours)	✓	FP16	LoRA	50.6
					NLS	52.5^{+1.9}
	SQFT (Ours)	✗	INT4 + FP16	LoRA	48.9	
				NLS	50.0^{+1.1}	
	SQFT + QA-SparsePEFT (Ours)	✓	INT4	LoRA	48.2	
				NLS	50.2^{+2.0}	
	Shears	✗	FP16 + FP16	LoRA	25.5	
				NLS	27.9^{+2.4}	
	SQFT + SparsePEFT (Ours)	✓	FP16	LoRA	22.1	
				NLS	24.9^{+2.8}	
	SQFT (Ours)	✗	INT4 + FP16	LoRA	24.2	
				NLS	25.2^{+1.0}	
	SQFT + QA-SparsePEFT (Ours)	✓	INT4	LoRA	22.6^{+0.2}	
				NLS	22.4	

made regarding model storage, inference speedup, or memory: merging is better than unmerging due to the overhead from the unmerged adapters, and quantization mode is better than full-precision mode. As for accuracy, the mergeable method we propose is competitive with the previous non-mergeable method. Regarding the fine-tuning time, our mergeable method is slightly slower than the non-mergeable method due to the additional mask and adapter calculations. In summary, SQFT with SparsePEFT is the best choice for full-precision mode because it eliminates the adapter’s additional path without sacrificing accuracy. Suppose memory usage during fine-tuning is a priority for the quantization mode. In that case, vanilla SQFT (first configuration in Figure 2) is the best choice because it only requires the quantized model with little overhead of different precision adapters. Otherwise, SQFT with QA-SparsePEFT is better because it can ultimately produce a most efficient model that will be of great benefit at deployment time.

3.6 Ablation Studies - LoRA vs NLS

As shown in Table 5, the ablation studies across 30%, 50%, and 70% sparsity highlight the benefits of elastic adapters (NLS), which enhance the performance of SQFT, further reducing the gap to the dense or non-quantized models while enjoying the advantages of sparsity or quantization.

Table 6: Cost analysis for different pipelines (**rank**). ID 1, 2, 3, and 4 represent LoRA/Shears, SQFT + SparsePEFT, SQFT, and SQFT + QA-SparsePEFT, respectively.

ID	1	2	3	4
Mergeable	✗	✓	✗	✓
Final Precision	FP16 + FP16	FP16	INT4 + FP16	INT4
Model Storage (↓)		1 > 2 > 3 > 4		
Fine-tuning Time (↓)		1 ≈ 3 < 2 ≈ 4		
Fine-tuning Memory (↓)		3 < 1 ≈ 2 ≈ 4		
Inference Speedup (↑)		4 > 3 > 2 > 1		
Inference Memory (↓)		4 < 3 < 2 < 1		
Accuracy (↑)		1 ≈ 2 > 3 ≈ 4		

4 Conclusion

Large pre-trained models often require fine-tuning to downstream target tasks and compression to utilize them in resource-constrained environments. This paper presents SQFT, a low-cost fine-tuning solution for low precision and sparse foundation models. SQFT solves challenges when merging sparse (and quantized) base models and dense (with different numerical precision) adapters without losing the induced sparsity in the base model while delivering high-performing fine-tuned models. We make a few SQFT’s fine-tuned models available to reviewers for reproducing our results at: https://anonymous.4open.science/r/sqft_examples-71C7

504 Limitations and Ethical Considerations

505 Large pre-trained models have gained popularity
506 and are the base of many applications. However,
507 these models are often used indiscriminately with
508 little analysis of their potential failures and conse-
509 quences. SQFT solely focuses on these large mod-
510 els’ efficient fine-tuning and compression. How-
511 ever, users of SQFT should also consider the lim-
512 itations of these models before deployment in en-
513 vironments where they can cause harm or conflict.
514 Although compressing and fine-tuning these mod-
515 els on a particular downstream task would make
516 them perform better, more studies are needed re-
517 garding the effects of this specialization.

518 We demonstrate SQFT on several pre-trained
519 models. The benefits obtained from the pro-
520 posed solution might transfer smoothly to other
521 transformer-based models. However, there might
522 also be models and datasets in which additional
523 considerations must be taken. For instance, in our
524 current experiments, we have noticed that in the
525 case of OpenELM-1.1B (Mehta et al., 2024), fine-
526 tuning on math reasoning datasets, e.g., GSM8K,
527 does not result in high accuracy, and more exper-
528 imentation is needed. There is also the case in
529 which a pre-trained model might have been trained
530 on a particular benchmark, a form of data contam-
531 ination, which is difficult to confirm since often
532 the details of the training data are not shared pub-
533 licly (Zhang et al., 2024). In these cases, inducing
534 sparsity might result in a drop in accuracy on that
535 particular benchmark.

536 Due to the many unknowns and complexity of
537 current large models, it is essential to take measures
538 to prevent their use in sensitive applications. With
539 insights obtained by the research community in
540 the years to come, understanding the intricacies of
541 these models will help us use them beneficially and
542 safely.

543 References

544 Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng
545 Gao, and Yejin Choi. 2020. Piqa: Reasoning about
546 physical commonsense in natural language. In *Thirty-
547 Fourth AAAI Conference on Artificial Intelligence*.

548 Christopher Clark, Kenton Lee, Ming-Wei Chang,
549 Tom Kwiatkowski, Michael Collins, and Kristina
550 Toutanova. 2019. Boolq: Exploring the surprising
551 difficulty of natural yes/no questions. In *NAACL*.

552 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,
553 Ashish Sabharwal, Carissa Schoenick, and Oyvind

Tafjord. 2018. [Think you have solved question an-
554 swering? try arc, the ai2 reasoning challenge](#). *ArXiv*,
555 abs/1803.05457. 556

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
557 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias
558 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro
559 Nakano, Christopher Hesse, and John Schulman.
560 2021. [Training verifiers to solve math word prob-
561 lems](#). *Preprint*, arXiv:2110.14168. 562

Tim Dettmers, Mike Lewis, Younes Belkada, and Luke
563 Zettlemoyer. 2022. [GPT3.int8\(\): 8-bit matrix mul-
564 tiplication for transformers at scale](#). In *Advances in
565 Neural Information Processing Systems*. 566

Tim Dettmers and Luke Zettlemoyer. 2023. The case
567 for 4-bit precision: k-bit inference scaling laws. In
568 *Proceedings of the 40th International Conference on
569 Machine Learning, ICML’23*. JMLR.org. 570

Elias Frantar and Dan Alistarh. 2023. SparseGPT: Mas-
571 sive language models can be accurately pruned in
572 one-shot. *arXiv preprint arXiv:2301.00774*. 573

Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and
574 Dan Alistarh. 2022a. [GPTQ: Accurate post-training
575 compression for generative pretrained transformers](#).
576 *arXiv preprint arXiv:2210.17323*. 577

Elias Frantar, Sidak Pal Singh, and Dan Alistarh. 2022b.
578 Optimal Brain Compression: a framework for ac-
579 curate post-training quantization and pruning. *Ad-
580 vances in Neural Information Processing Systems*,
581 36. 582

Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman,
583 Sid Black, Anthony DiPofi, Charles Foster, Laurence
584 Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li,
585 Kyle McDonell, Niklas Muennighoff, Chris Ociepa,
586 Jason Phang, Laria Reynolds, Hailey Schoelkopf,
587 Aviya Skowron, Lintang Sutawika, Eric Tang, An-
588 ish Thite, Ben Wang, Kevin Wang, and Andy Zou.
589 2023. [A framework for few-shot language model
590 evaluation](#). 591

Masafumi Hagiwara. 1994. [A simple and effective
592 method for removal of hidden units and weights](#). *Neu-
593 rocomputing*, 6(2):207–218. Backpropagation, Part
594 IV. 595

Torsten Hoefler, Dan Alistarh, Tal Ben-Nun, Nikoli
596 Dryden, and Alexandra Peste. 2021. Sparsity in deep
597 learning: pruning and growth for efficient inference
598 and training in neural networks. *J. Mach. Learn. Res.*,
599 22(1). 600

Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-
601 Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu
602 Chen. 2022. [LoRA: Low-rank adaptation of large
603 language models](#). In *International Conference on
604 Learning Representations*. 605

Zhiqiang Hu, Yihuai Lan, Lei Wang, Wanyu Xu, Ee-
606 Peng Lim, Roy Ka-Wei Lee, Lidong Bing, and Sou-
607 janya Poria. 2023. Llm-adapters: An adapter family
608

609	for parameter-efficient fine-tuning of large language models. <i>arXiv preprint arXiv:2304.01933</i> .	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	666
610		Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	667
611	Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate	Kaiser, and Illia Polosukhin. 2017. Attention is all	668
612	Kushman, and Hannaneh Hajishirzi. 2016. MAWPS: A math word problem repository . In <i>Proceedings of</i>	<i>you need</i> . In <i>Advances in Neural Information Pro-</i>	669
613	<i>the 2016 Conference of the North American Chapter</i>	<i>cessing Systems</i> , volume 30. Curran Associates, Inc.	670
614	<i>of the Association for Computational Linguistics: Hu-</i>		
615	<i>man Language Technologies</i> , pages 1152–1157, San	Peisong Wang, Qiang Chen, Xiangyu He, and Jian	671
616	Diego, California. Association for Computational	Cheng. 2020. Towards accurate post-training net-	672
617	Linguistics.	work quantization via bit-split and stitching . In <i>Pro-</i>	673
618		<i>ceedings of the 37th International Conference on Ma-</i>	674
619	Yann LeCun, John Denker, and Sara Solla. 1989. Op-	<i>chine Learning</i> , volume 119 of <i>Proceedings of Ma-</i>	675
620	timal brain damage . In <i>Advances in Neural In-</i>	<i>chine Learning Research</i> , pages 9847–9856. PMLR.	676
621	<i>formation Processing Systems</i> , volume 2. Morgan-		
622	Kaufmann.	Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu,	677
623	Sachin Mehta, Mohammad Sekhvat, Qingqing Cao,	Julien Demouth, and Song Han. 2023. SmoothQuant:	678
624	Max Horton, Yanzi Jin, Frank Sun, Iman Mirzadeh,	Accurate and efficient post-training quantization for	679
625	Mahyar Najibikohnehshahri, Dmitry Belenko, Pe-	large language models . In <i>Proceedings of the 40th</i>	680
626	ter Zatloukal, and Mohammad Rastegari. 2024. Openelm:	<i>International Conference on Machine Learning</i> .	681
627	An efficient language model family with		
628	open training and inference framework .	Peng Xu, Wenqi Shao, Mengzhao Chen, Shitao Tang,	682
629	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish	Kaipeng Zhang, Peng Gao, Fengwei An, Yu Qiao,	683
630	Sabharwal. 2018. Can a suit of armor conduct elec-	and Ping Luo. 2024. Besa: Pruning large language	684
631	tricity? a new dataset for open book question answer-	models with blockwise parameter-efficient sparsity	685
632	ing . In <i>Conference on Empirical Methods in Natural</i>	allocation . <i>Preprint</i> , arXiv:2402.16880.	686
633	<i>Language Processing</i> .	Zhewei Yao, Reza Yazdani Aminabadi, Minjia Zhang,	687
634	J. Pablo Munoz, Jinjie Yuan, and Nilesh Jain. 2024a. Shears: Unstructured sparsity with neural low-rank	Xiaoxia Wu, Conglong Li, and Yuxiong He. 2022. Zeroquant: Efficient and affordable post-training	688
635	adapter search . <i>The 2024 Annual Conference of the</i>	quantization for large-scale transformers . In <i>Ad-</i>	689
636	<i>North American Chapter of the Association for Com-</i>	<i>vances in Neural Information Processing Systems</i> ,	690
637	<i>putational Linguistics (NAACL-2024)</i> .	volume 35, pages 27168–27183. Curran Associates,	691
638		Inc.	692
639	J. Pablo Munoz, Jinjie Yuan, Yi Zheng, and Nilesh Jain.	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali	694
640	2024b. LoNAS: Elastic low-rank adapters for effi-	Farhadi, and Yejin Choi. 2019. Hellaswag: Can a	695
641	cient large language models . In <i>Proceedings of the</i>	machine really finish your sentence? In <i>Proceedings</i>	696
642	<i>the 2024 Joint International Conference on Computa-</i>	<i>of the 57th Annual Meeting of the Association for</i>	697
643	<i>tional Linguistics, Language Resources and Evalu-</i>	<i>Computational Linguistics</i> .	698
644	<i>ation (LREC-COLING 2024)</i> , pages 10760–10776,	Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson,	699
645	Torino, Italia. ELRA and ICCL.	Catherine Wu, Will Song, Tiffany Zhao, Pranav Raja,	700
646	Markus Nagel, Rana Ali Amjad, Mart Van Baalen,	Dylan Slack, Qin Lyu, Sean Hendryx, Russell Kap-	701
647	Christos Louizos, and Tijmen Blankevoort. 2020. Up	lan, Michele Lunati, and Summer Yue. 2024. A	702
648	or down? adaptive rounding for post-training quanti-	careful examination of large language model per-	703
649	zation . In <i>Proceedings of the 37th International Con-</i>	formance on grade school arithmetic . <i>Preprint</i> ,	704
650	<i>ference on Machine Learning, ICML’20</i> . JMLR.org.	arXiv:2405.00332.	705
651	Arkil Patel, Satwik Bhattamishra, and Navin Goyal.	Mingyang Zhang, Hao Chen, Chunhua Shen, Zhen	706
652	2021. Are NLP models really able to solve simple	Yang, Linlin Ou, Xinyi Yu, and Bohan Zhuang.	707
653	math word problems? In <i>Proceedings of the 2021</i>	2023. Loraprune: Pruning meets low-rank parameter-	708
654	<i>Conference of the North American Chapter of the</i>	efficient fine-tuning . <i>Preprint</i> , arXiv:2305.18403.	709
655	<i>Association for Computational Linguistics: Human</i>		
656	<i>Language Technologies</i> , pages 2080–2094, Online.		
657	Association for Computational Linguistics.		
658	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavat-		
659	ula, and Yejin Choi. 2021. Winogrande: An adver-		
660	sarial winograd schema challenge at scale . <i>Commun.</i>		
661	<i>ACM</i> , 64(9):99–106.		
662	Mingjie Sun, Zhuang Liu, Anna Bair, and J. Zico		
663	Kolter. 2023. A simple and effective pruning ap-		
664	proach for large language models . <i>arXiv preprint</i>		
665	<i>arXiv:2306.11695</i> .		

Appendix

A Related Work

Generative pre-trained models often based on the Transformer architecture (Vaswani et al., 2017) require the application of compression techniques to reduce their significant computational cost and to address challenges, e.g., related to memory bandwidth. Classic compression techniques like pruning and quantization have been adapted to the age of LPMs, removing inefficiencies that cannot be tolerated when dealing with billions of parameters. We discuss them in more detail next.

Pruning Inducing sparsity, either by zeroing out weights or activations or removing network elements, can improve the efficiency of LPMs during inference, provided that they are executed on a runtime that can exploit sparse patterns. Pruning has a long history (LeCun et al., 1989), but with the advent of LPMs, traditional methods (Hoefler et al., 2021), e.g., Magnitude Pruning (Hagiwara, 1994), have been replaced by new approaches that are suited for the challenges of these models. In particular, due to their large number of parameters. SparseGPT (Frantar and Alistarh, 2023) proposes a one-shot pruning method for transformer-based models that trade minimal accuracy drop for increasing sparsity levels. The method approaches LPMs’ pruning layer-wise with an efficient weight reconstruction algorithm that incrementally prunes the weight matrix elements. Wanda (Sun et al., 2023) proposes a more straightforward approach that does not require weight updates, computing a score using the weight magnitude and the norm of input activations. This approach obtains better results than SparseGPT. Recently, BESA (Xu et al., 2024) improves over SparseGPT and Wanda by targeting individual transformer blocks and allocating sparsity per layer using a differentiable method. These approaches induce sparsity on pre-trained models and are evaluated on zero-shot benchmarks. Our end-to-end solution, SQFT, focuses on further adapting the sparsified models to new tasks or datasets.

Quantization In the era of large pre-trained foundation/frontier models (LPMs), quantization approaches have evolved to address the challenges of scale and memory bandwidth. Due to the high cost of retraining these models to recover accuracy degradation, special consideration has to

be taken when incorporating compression techniques, like quantization-aware training in foundation models. Post-training, one-shot quantization methods have prevailed, obtaining quantized versions of large models in hours. LLM.Int8() was among the first Int8 quantization procedures for large-scale transformer-based PLMs (Dettmers et al., 2022). Using vector-wise quantization and mixed-precision decomposition, LLM.Int8() demonstrated that it can effectively confront the outliers that emerge in activations, which makes traditional quantization methods fail in models with more than 6.7B parameters. In a contemporary work, after running thousands of experiments with various large pre-trained models, it was demonstrated that 4-bit parameters can reach optimal performance compared to other bit-precisions in the 3 to 16-bit range (Dettmers and Zettlemoyer, 2023). ZeroQuant (Yao et al., 2022) quantizes GPT-3 models, obtaining a reduction in latency up to 4.16x by utilizing group-wise quantization for weights, token-wise quantization for activations, and layer-by-layer knowledge distillation. SmoothQuant (Xiao et al., 2023) makes activations easier to quantize by smoothing them and compensating this operation with a transformation of the weights, resulting in improved results over ZeroQuant and LLM.Int8(). GPTQ is another good representative of one-shot quantization approaches designed especially for LPMs (Frantar et al., 2022a). GPTQ builds on the learnings from Optimal Brain Quantization (OBQ) (Frantar et al., 2022b) and applies layer-wise quantization to the full-precision weights of a base LPM. We incorporate GPTQ as the default quantization method in SQFT’s pre-fine-tuning stage.

Parameter-efficient Fine-tuning (PEFT) Due to their large number of parameters, it is too costly to fine-tune pre-trained large models. Updating all their weights to improve their performance in a downstream task might require devices with large memory capacity. PEFT techniques attempt to address this challenge by avoiding the update of all weights in the pre-trained model. For instance, low-rank (LoRA) adapters (Hu et al., 2022) use a fraction (often less than 1%) of additional weights to adapt the model to a new task. LoRA adapters, B and A , are utilized to reparameterize a linear projection, $Y = WX$, keeping the weights, W , frozen and updating only the low-rank adapter matrices, A and B , i.e., $Y = WX + BAX$.

Algorithm 1 Hill-climbing Search Algorithm

Input: Number of turns T , Number of neighbors N , Neighbor step size S , Number of evaluation samples M , Heuristic configuration c_h , Validation dataset \mathcal{D}

Output: Optimal configuration c^*

```
1:  $c_a \leftarrow c_h$   $\triangleright$  Initialize anchor with the heuristic configuration
2:  $V \leftarrow \{c_h\}$   $\triangleright$  Initialize the set of visited configurations
3:  $\mathcal{D}_M \leftarrow \text{Sample}(\mathcal{D}, M)$   $\triangleright$  Create a proxy dataset by randomly sampling  $M$  samples from  $\mathcal{D}$ 
4: for  $t \leftarrow 1$  to  $T$  do
5:    $\mathcal{C} \leftarrow \text{Neighbor-sample}(c_a, N, S) - V$   $\triangleright$  Sample  $N$  unvisited  $S$ -step neighbor configs
6:    $\mathcal{V} \leftarrow V \cup \mathcal{C}$   $\triangleright$  Add the sampled configurations to the set of visited configurations
7:    $c_m \leftarrow \text{MaxAcc}(\text{Eval}(\mathcal{D}_M, \mathcal{C}))$   $\triangleright$  The config with the maximum accuracy on proxy data
8:   if  $\text{Acc}(c_m) > \text{Acc}(c^*)$  then
9:      $c_a \leftarrow c_m$   $\triangleright$  Update anchor configuration if the new configuration has higher accuracy
10:  end if
11: end for
12:  $c^* \leftarrow c_a$   $\triangleright$  The optimal configuration is the final anchor configuration
13: return  $c^*$ 
```

810 Recently, Shears proposed Neural Low-rank
811 Adapter Search (Munoz et al., 2024a) and demon-
812 strated that LoRA adapters can be made elastic to
813 allow for the application of weight-sharing schemes
814 and keeping the original weights of the model
815 frozen and compressed, e.g., inducing sparsity be-
816 fore the fine-tuning stage. However, a challenge
817 that emerges is that merging the dense adapters
818 with the sparse weights results in the overall loss of
819 sparsity. LoRAPrune has attempted to address this
820 challenge by using the weights and gradients of the
821 LoRA adapters to remove elements in the model’s
822 weights (Zhang et al., 2023). As demonstrated in
823 the main sections of the paper, SQFT proposes an
824 alternative method for merging the dense adapters
825 with a minimal drop in accuracy.

shows that up to high sparsity levels, SQFT delivers
high-performing models.

842
843

B Hyperparameters

826
827 The hyperparameters used in our main experiments
828 are shown in Table 7.

C Hill-climbing search algorithm

829
830 We propose Algorithm 1 to start from the refer-
831 ence configuration (Section 3.1) and systematically
832 explore its neighbors. Table 4 in the main pa-
833 per shows the benefits of using any available bud-
834 get to execute this algorithm and discover better-
835 performing models.

D Additional Sparsity Levels and Ablation Studies for Llama-3 on GSM8K

836
837
838
839 We conducted additional experiments and ablations
840 studies with different sparsity levels and compared
841 the underlying NLS approach to LoRA. Table 8

Table 7: Hyperparameters used in our experiments. For all approaches with NLS, we explored several manually designed search spaces and identified the optimal configuration for each pipeline. Note that in our experiments involving GSM8K and math instruction tuning, we conducted trials over 3 or 4 epochs and reported the best results achieved. Interestingly, SQFT with QA-SparsePEFT often necessitates extended training periods to exploit its quantization-aware capabilities fully.

Model	Task	Sparsity	Method	Epoch	Batch size	Learning rate	Adapter rank	Adapter alpha	Adapter target modules
Llama-3-8B	GSM8K	50%	LoRA	3	16	3e-4	32	64	Q, K, V, Up, Down
			Shears	3	16	3e-4	32,28,24,20,16	64	Q, K, V, Up, Down
			SQFT + SparsePEFT	3	16	3e-4	48,32,16	64	Q, K, V, Up, Down
			GPTQ + LoRA	3	16	3e-4	32	64	Q, K, V, Up, Down
			SQFT	3	16	3e-4	40,32,24	64	Q, K, V, Up, Down
			SQFT + QA-SparsePEFT	4	16	3e-4	48,32,16	64	Q, K, V, Up, Down
Phi-3-Mini-4K-Instruct	Math	50%	LoRA	3	16	3e-4	32	64	Qkv
			Shears	3	16	3e-4	48,40,32,24,16	64	Qkv
			SQFT + SparsePEFT	3	16	3e-4	48,32,16	64	Qkv
			GPTQ + LoRA	3	16	3e-4	32	64	Qkv
			SQFT	3	16	3e-4	32,28,24,20,16	64	Qkv
			SQFT + QA-SparsePEFT	4	16	3e-4	32,24,16	64	Qkv
Phi-3-Mini-4K-Instruct	CS	50%	LoRA	3	16	1e-4	16	32	Qkv
			Shears	3	16	1e-4	16,12,8	32	Qkv
			SQFT + SparsePEFT	3	16	1e-4	16,12,8	32	Qkv
			GPTQ + LoRA	3	16	1e-4	16	32	Qkv
			SQFT	3	16	1e-4	16,12,8	32	Qkv
			SQFT + QA-SparsePEFT	3	16	1e-4	16,12,8	32	Qkv

Table 8: Ablation studies for various sparsity levels (**Llama-3-8B** with GSM8K).

Model	Sparsity	Method	Mergeable	Final Precision (Base + Adapter / Base)	Fine-tune Approach	GSM8K Test Accuracy(%)
Llama-3-8B	0%	w/o tune	-	FP16	-	50.0
		w/o tune	-	<i>w/o Quantization</i> FP16	-	47.5
	20%	Shears	✗	FP16 + FP16	LoRA NLS	58.7 61.2^{+2.5}
		SQFT + SparsePEFT	✓	FP16	LoRA NLS	60.3 62.0^{+1.7}
		w/o tune	-	<i>Quantization</i> INT4	-	36.6
		SQFT	✗	INT4 + FP16	LoRA NLS	57.8 60.0^{+2.2}
		SQFT + QA-SparsePEFT	✓	INT4	LoRA NLS	54.7 55.6^{+0.9}
		w/o tune	-	<i>w/o Quantization</i> FP16	-	40.9
	30%	Shears	✗	FP16 + FP16	LoRA NLS	58.2 59.8^{+1.6}
		SQFT + SparsePEFT	✓	FP16	LoRA NLS	60.0 61.2^{+1.2}
		w/o tune	-	<i>Quantization</i> INT4	-	30.3
		SQFT	✗	INT4 + FP16	LoRA NLS	56.7 57.6^{+0.9}
		SQFT + QA-SparsePEFT	✓	INT4	LoRA NLS	54.8 56.0^{+1.2}
		w/o tune	-	<i>w/o Quantization</i> FP16	-	31.6
	40%	Shears	✗	FP16 + FP16	LoRA NLS	56.9 57.9^{+1.5}
		SQFT + SparsePEFT	✓	FP16	LoRA NLS	56.9 56.4
		w/o tune	-	<i>Quantization</i> INT4	-	20.1
		SQFT	✗	INT4 + FP16	LoRA NLS	54.9 54.9
		SQFT + QA-SparsePEFT	✓	INT4	LoRA NLS	53.4 53.7^{+0.3}
		w/o tune	-	<i>w/o Quantization</i> FP16	-	12.5
	50%	Shears	✗	FP16 + FP16	LoRA NLS	50.6 52.2^{+1.6}
		SQFT + SparsePEFT	✓	FP16	LoRA NLS	50.6 52.5^{+1.9}
		w/o tune	-	<i>Quantization</i> INT4	-	7.0
		SQFT	✗	INT4 + FP16	LoRA NLS	48.9 50.0^{+1.1}
		SQFT + QA-SparsePEFT	✓	INT4	LoRA NLS	48.2 50.2^{+2.0}
		w/o tune	-	<i>w/o Quantization</i> FP16	-	39.9
	60%	Shears	✗	FP16 + FP16	LoRA NLS	39.9 45.3^{+5.4}
		SQFT + SparsePEFT	✓	FP16	LoRA NLS	40.7 42.5^{+1.8}
		w/o tune	-	<i>Quantization</i> INT4	-	40.1
		SQFT	✗	INT4 + FP16	LoRA NLS	48.9 42.0^{+1.9}
		SQFT + QA-SparsePEFT	✓	INT4	LoRA NLS	37.6 40.9^{+3.3}
		w/o tune	-	<i>w/o Quantization</i> FP16	-	25.5
	70%	Shears	✗	FP16 + FP16	LoRA NLS	25.5 27.9^{+2.4}
		SQFT + SparsePEFT	✓	FP16	LoRA NLS	22.1 24.9^{+2.8}
		w/o tune	-	<i>Quantization</i> INT4	-	24.2
		SQFT	✗	INT4 + FP16	LoRA NLS	24.2 25.2^{+1.0}
		SQFT + QA-SparsePEFT	✓	INT4	LoRA NLS	22.4 22.0^{+0.2}
		w/o tune	-	<i>w/o Quantization</i> FP16	-	22.4