# LittleBit: Ultra Low-Bit Quantization via Latent Factorization

**Banseok Lee**[*]  **Dongkyu Kim**[*]  **Youngcheon You**  **Youngmin Kim**[†]
Samsung Research
{bs93.lee, dongkyu.k, y01000.you, ym1012.kim}@samsung.com

## Abstract

Deploying large language models (LLMs) often faces challenges from substantial memory and computational costs. Quantization offers a solution, yet performance degradation in the sub-1-bit regime remains particularly difficult. This paper introduces LittleBit, a novel method for extreme LLM compression. It targets levels like 0.1 bits per weight (BPW), achieving nearly *31× memory reduction*, *e.g.*, Llama2-13B to under 0.9 GB. LittleBit represents weights in a low-rank form using latent matrix factorization, subsequently binarizing these factors. To counteract information loss from this extreme precision, it integrates a multi-scale compensation mechanism. This includes row, column, and an *additional* latent dimension that learns per-rank importance. Two key contributions enable effective training: Dual Sign-Value-Independent Decomposition (Dual-SVID) for quantization-aware training (QAT) initialization, and integrated Residual Compensation to mitigate errors. Extensive experiments confirm LittleBit's superiority in sub-1-bit quantization: *e.g.*, its 0.1 BPW performance on Llama2-7B surpasses the leading method's 0.7 BPW. LittleBit establishes a new, viable size-performance trade-off–unlocking a potential *11.6× speedup* over FP16 at the kernel level–and makes powerful LLMs practical for resource-constrained environments.

## 1  Introduction

Large Language Models (LLMs) based on the Transformer architecture [1] have revolutionized natural language processing. However, their immense scale, reaching hundreds of billions or trillions of parameters [2, 3], leads to prohibitive computational and memory costs, particularly GPU VRAM, hindering their widespread deployment on consumer or edge devices [4–6].

Model quantization, which reduces numerical precision, is a key technique to address these bottlenecks [7]. While post-training quantization (PTQ) methods like GPTQ [8] and AWQ [9] compress to ∼4 bits, further compression (*e.g.*, to 1-bit) typically requires quantization-aware training (QAT) [10, 11] to maintain performance. Indeed, QAT has demonstrated the ability to enable substantial 1-bit compression while preserving model fidelity [12–14].

However, even 1-bit models (*e.g.*, ∼15.4GB for 70B parameters [7]) can exceed the capacity of highly resource-constrained devices [15], motivating the exploration of sub-1-bit quantization. Although prior work [16] has ventured into this area, maintaining performance at extremely low effective bits (*e.g.*, 0.55 bits per weight (BPW)) while ensuring computational efficiency without restrictive hardware dependencies remains a critical challenge [16]. Consequently, it is crucial to pioneer new techniques that can achieve such profound model compression, aiming for bit rates around 0.1 BPW, while maintaining the integrity of model performance. Our approach stems from two

---

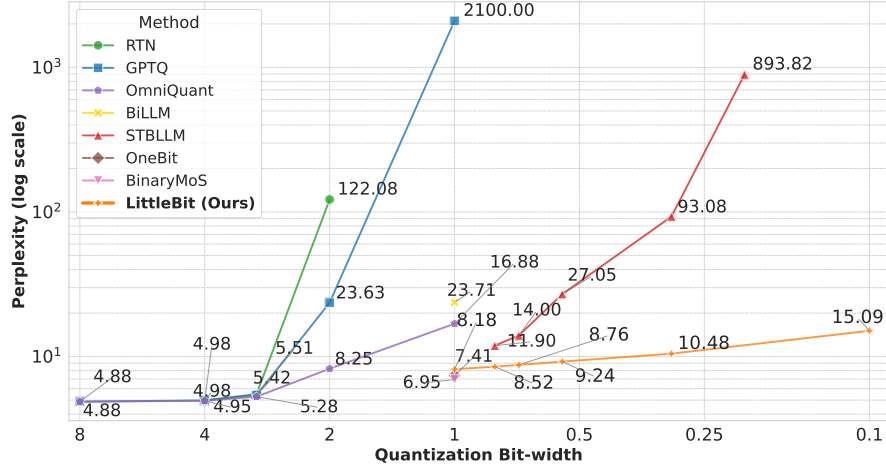[*]Equal contribution
[†]Corresponding author

Figure 1: Low-bit quantization perplexity on Llama2-13B (WikiText-2). LittleBit surpasses the state-of-the-art sub-1-bit quantization technique. Below 0.5 BPW, where the leading prior method degrades sharply, ours remains robust down to 0.1 BPW.

observations. First, LLM weight matrices often exhibit significant low-rank structure [17, 18]. This suggests that factorization methods like Singular Value Decomposition (SVD) [19] could offer a more stable compression pathway than pruning [20], especially under extreme compression regimes. For instance, at compression ratios exceeding 50%, SVD-based approaches [21] consistently outperform pruning methods like Wanda [22] (see Section A.5 for details). Second, binarization inherently causes information loss [23]. Recent high-performing 1-bit methods [13, 24, 25] demonstrate that sophisticated scaling over multiple dimensions (*e.g.*, row and column), is crucial for mitigating this loss and stabilizing performance [25].

Based on these insights, we introduce **LittleBit**, a novel method for extreme sub-1-bit quantization (*e.g.*, 0.1 BPW). LittleBit first represents weights via low-rank factorization ($\mathbf{W} \approx \mathbf{UV}^\top$) and then binarizes these factors. To counteract errors from binarization, LittleBit leverages a multi-scale compensation mechanism that applies learnable scales across rows, columns, and an *additional latent* dimension. This entire architecture is then complemented by Residual Compensation.

Extensive experiments show LittleBit's superiority over STBLLM [16], the leading sub-1-bit technique, on LLMs from 1.3B to 32B parameters. Notably, LittleBit achieves competitive performance on Llama2-13B at an unprecedented 0.1 BPW (Fig. 1), surpassing STBLLM at 0.55 BPW. Furthermore, the 32B LittleBit model at just 0.3 BPW maintains strong performance, significantly outperforming STBLLM and setting a new state-of-the-art for sub-1-bit quantization.

In summary, our contribution is as follows:

- We propose LittleBit, a novel method unifying latent matrix factorization with multi-scale compensation for extreme sub-1-bit quantization. Our method achieves unprecedentedly low effective bits (*e.g.*, down to 0.1 BPW) while preserving model performance, establishing a new, viable size-performance trade-off that provides a practical path for deploying LLMs in highly resource-constrained environments.
- We introduce Dual-SVID for effective initialization of factorized structures and Residual Compensation to mitigate errors in extreme quantization.
- We provide extensive validation showing that LittleBit significantly outperforms STBLLM [16], the prior leading sub-1-bit quantization technique, on various LLM scales.

## 2 Related Works

### 2.1 Binarization and Quantization

Network binarization targets extreme compression and acceleration by constraining weights or activations to $\pm 1$, enabling efficient bitwise operations [26, 27]. Although early works like BinaryConnect [28] and BNNs [26] showed feasibility, they faced significant accuracy degradation [29].
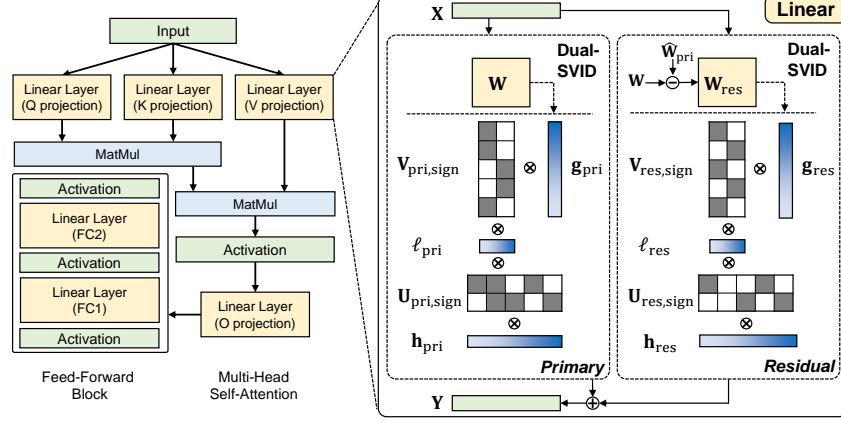
Figure 2: Comparison of a standard Transformer layer (left) and the LittleBit architecture (right). LittleBit performs linear transformation using parallel Primary and Residual pathways. The Primary path employs binarized factors ($\mathbf{U}_{\text{sign}}, \mathbf{V}_{\text{sign}}$) and FP16 scales ($\mathbf{h}, \mathbf{g}, \boldsymbol{\ell}$) on input $\mathbf{X}$, initialized from $\mathbf{W}$ via Dual-SVID. Simultaneously, the Residual path computes a correction with its own parameters ($\mathbf{U}_{\text{res,sign}}, \mathbf{V}_{\text{res,sign}}, \mathbf{h}_{\text{res}}, \mathbf{g}_{\text{res}}, \boldsymbol{\ell}_{\text{res}}$) from the approximation residual. Their outputs sum to form $\mathbf{Y}$, eliminating storage of the effective weight matrices $\widehat{\mathbf{W}}_{\text{pri}}$ and $\widehat{\mathbf{W}}_{\text{res}}$.

This degradation was partially mitigated by introducing scaling factors [30] and employing the Straight-Through Estimator (STE) [31] to handle the non-differentiable quantization function during training [32].

However, directly applying these early techniques to sensitive large language models (LLMs) typically results in unacceptable performance loss [33]. To address this challenge, LLM-specific quantization strategies have been developed, mainly categorized into post-training quantization (PTQ) and quantization-aware training (QAT) [34]. PTQ adapts pre-trained models, achieving competitive results around 4-bit precision with methods such as GPTQ [8] and AWQ [9], but generally struggles to maintain performance below 2-bit precision [33, 7]. The sub-1-bit regime is particularly challenging for PTQ methods due to severe information loss [35]. Recent work such as STBLLM [16] shows that incorporating structured binarized compression can push PTQ into the sub-1-bit regime, yet a non-negligible accuracy gap remains.

In contrast, QAT integrates the quantization process directly into the training loop. This enables the model to learn and adapt to the low-precision format during training, generally yielding superior performance compared to PTQ at very low bit-widths [36]. Successful QAT approaches for low-bit scenarios often employ sophisticated scaling mechanisms alongside the adaptive training itself, as evidenced by several recent methods [13, 24]. However, despite these advancements, consistently achieving high model fidelity when quantizing to less than 1.0 BPW remains a considerable challenge. This persisting difficulty underscores the need for novel architectures and training strategies specifically designed for such extreme compression regimes.

## 2.2 Low-Rank Approximation in Quantization

Beyond quantization, the inherent parameter redundancy within deep neural networks [35], particularly evident in LLMs [37, 38], can be effectively exploited using low-rank approximation methods. Techniques such as Singular Value Decomposition (SVD) [19] allow large weight matrices to be represented as products of smaller factors, offering an alternative compression strategy that can be more resilient than pruning, especially under high compression regimes [20]. Recognizing their potential for complementary benefits, researchers have explored combining low-rank factorization to reduce parameter count and quantization to lower parameter precision [39]. Initial studies applied SVD in various roles within the LLM context: some used activation statistics to guide decomposition [21], others integrated SVD principles into initialization or parameter-efficient updates [40, 41, 13], applied transformations prior to factorization [42], or even employed low-rank structures to model quantization errors [43]. A more integrated approach involves incorporating low-rank concepts directly within the QAT methods, as demonstrated by DL-QAT [44], which jointly learns a representation amenable to both low-rank approximation and low-precision quantization, albeit limited to 3-bit

quantization. This direction of optimizing structural properties such as rank and numerical precision during training represents a promising path towards more effective LLM compression, with potential for further exploration of even lower precision quantization.

# 3 Methodology

This section details **LittleBit**. We describe its factorized architecture with multi-scale compensation (Section 3.1), the Dual-SVID initialization for stable QAT (Section 3.2), and Residual Compensation for enhanced fidelity (Section 3.3). These components are optimized via QAT (details in Section 4).

## 3.1 LittleBit Architecture

The LittleBit architecture, depicted in Fig. 2, redesigns linear layers for extreme compression (*e.g.*, 0.1 BPW) by synergistically combining low-rank factorization and binarization. LittleBit first leverages the observed low-rank structure in LLM weights [45] to approximate the weight matrix $\mathbf{W} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ as:

$$\mathbf{W} \approx \mathbf{U}\mathbf{V}^{\top} \quad (\mathbf{U} \in \mathbb{R}^{d_{\text{out}} \times r}, \mathbf{V} \in \mathbb{R}^{d_{\text{in}} \times r}, r \ll \min(d_{\text{in}}, d_{\text{out}})) \tag{1}$$

This factorization-based approach is chosen over other structural compression methods, such as pruning, due to its superior robustness in extreme compression regimes. As demonstrated in our comparative analysis against methods like Wanda (Section A.5), SVD-based compression maintains model integrity far more effectively at high compression ratios, providing a more stable foundation for our ultra low-bit quantization. The resulting factors $\mathbf{U}$ and $\mathbf{V}$ are then binarized:

$$\mathbf{U}_{\text{sign}} = \text{sign}(\mathbf{U}), \quad \mathbf{V}_{\text{sign}} = \text{sign}(\mathbf{V}) \in \{-1, +1\} \tag{2}$$

To compensate for the significant information loss incurred by binarization [23], LittleBit incorporates learnable FP16 scales. Beyond standard row ($\mathbf{h}$) and column ($\mathbf{g}$) scaling [33, 13, 25], it introduces an *additional latent* scale ($\boldsymbol{\ell} \in \mathbb{R}^r$). This latent scale is designed to address the factorization structure by learning the relative importance of each of the $r$ latent dimensions, which correspond to the columns of the binarized factors $\mathbf{U}_{\text{sign}}$ and $\mathbf{V}_{\text{sign}}$:

$$\mathbf{h} \in \mathbb{R}^{d_{\text{out}}}, \quad \mathbf{g} \in \mathbb{R}^{d_{\text{in}}}, \quad \boldsymbol{\ell} \in \mathbb{R}^r \tag{3}$$

The primary effective weight, $\widehat{\mathbf{W}}_{\text{pri}}$, is implicitly constructed from these components:

$$\widehat{\mathbf{W}}_{\text{pri}} = \text{diag}(\mathbf{h})\mathbf{U}_{\text{sign}}\text{diag}(\boldsymbol{\ell})\mathbf{V}_{\text{sign}}^{\top}\text{diag}(\mathbf{g}) \tag{4}$$

Instead of storing or optimizing $\widehat{\mathbf{W}}_{\text{pri}}$ directly, LittleBit learns latent full-precision factors $\mathbf{U}$ and $\mathbf{V}$, which are binarized during the forward pass using Eq. (2), along with the FP16 scales $\mathbf{h}, \mathbf{g}, \boldsymbol{\ell}$. This factorized representation enables efficient computation during the forward pass.

**Proposition 1** Let $\mathbf{X} \in \mathbb{R}^{\text{seq} \times d_{\text{in}}}$ be the input matrix. The forward computation $\mathbf{Y} = \mathbf{X}\widehat{\mathbf{W}}_{\text{pri}}^{\top}$ using the primary approximation (Eq. (4)) can be efficiently computed as:

$$\mathbf{Y} = ((((\mathbf{X} \odot \mathbf{g})\mathbf{V}_{\text{sign}}) \odot \boldsymbol{\ell})\mathbf{U}_{\text{sign}}^{\top}) \odot \mathbf{h} \tag{5}$$

where $\odot$ denotes element-wise multiplication with broadcasting.

This decomposition (Eq. (5)) replaces a large high precision General Matrix Multiply (GEMM) operation with two smaller binary matrix multiplications and element-wise scaling operations, offering computational advantages and reducing memory requirements. The effective bits per weight (BPW) is determined by the storage cost of these learnable parameters relative to the original matrix size (see Section A.4 for more details).

## 3.2 Dual-SVID Initialization

Naively initialing the highly constrained LittleBit structure can lead to unstable QAT. To circumvent this, we propose **Dual-SVID**, an SVD-based initialization method designed to provide a starting point where the initial effective weight $\widehat{\mathbf{W}}_{\text{pri},0}$ closely approximates $\mathbf{W}$. Dual-SVID aims to preserve essential sign and magnitude information from the optimal low-rank SVD factors (obtained
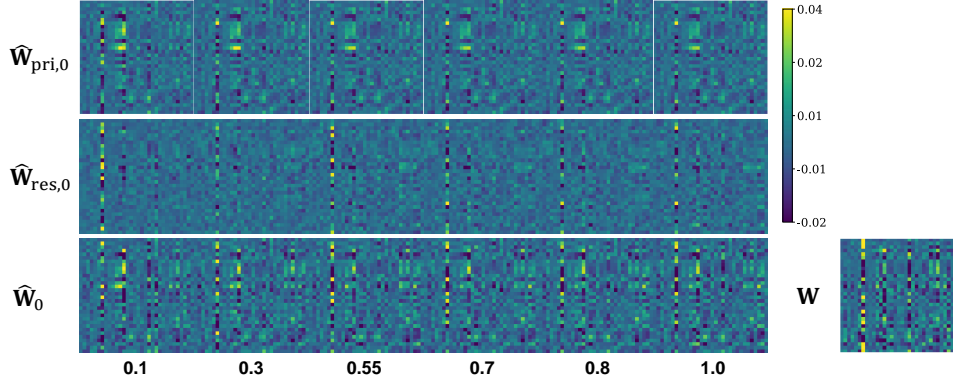
Figure 3: Visualization of Dual-SVID initialized weight components for a selected layer in Llama2-7B (Query Weight, Layer 0). Columns, from left to right, represent effective bits of 0.1, 0.3, 0.55, 0.7, 0.8, and 1.0 BPW. Rows display the primary approximation ($\widehat{\mathbf{W}}_{\text{pri},0}$), the residual approximation ($\widehat{\mathbf{W}}_{\text{res},0}$), and their sum ($\widehat{\mathbf{W}}_0 = \widehat{\mathbf{W}}_{\text{pri},0} + \widehat{\mathbf{W}}_{\text{res},0}$). The rightmost image shows the corresponding crop of the original weight matrix ($\mathbf{W}$) for reference.

from $\mathbf{W} \approx \mathbf{U}'\mathbf{V}'^{\top}$, then truncated) and map this information to the learnable LittleBit parameters ($\mathbf{U}_{\text{sign}}, \mathbf{V}_{\text{sign}}, \mathbf{h}, \mathbf{g}, \boldsymbol{\ell}$). The process involves: (1) Obtaining $\mathbf{U}', \mathbf{V}'$ via a truncated SVD. (2) Setting initial binary factors based on signs (Eq. (6)). (3) Decomposing the magnitudes $|\mathbf{U}'| \in \mathbb{R}^{d_{\text{out}} \times r}$ and $|\mathbf{V}'| \in \mathbb{R}^{d_{\text{in}} \times r}$. To initialize the scales, we perform a rank-1 SVD (or equivalent rank-1 approximation) on each of these magnitude matrices separately. For $|\mathbf{U}'|$, its best rank-1 approximation can be written as $\mathbf{s}_{U,0}(\boldsymbol{\ell}_{u,0})^{\top}$, where $\mathbf{s}_{U,0} \in \mathbb{R}^{d_{\text{out}}}$ becomes the initial row scale $\mathbf{h}_0$, and $\boldsymbol{\ell}_{u,0} \in \mathbb{R}^r$ is a component contributing to the latent scale. Similarly, for $|\mathbf{V}'|$, its rank-1 approximation $\mathbf{s}_{V,0}(\boldsymbol{\ell}_{v,0})^{\top}$ provides the initial column scale $\mathbf{g}_0 = \mathbf{s}_{V,0} \in \mathbb{R}^{d_{\text{in}}}$ and another latent scale component $\boldsymbol{\ell}_{v,0} \in \mathbb{R}^r$. These steps effectively disentangle multi-dimensional magnitude variations and allow us to estimate initial row ($\mathbf{h}_0$), column ($\mathbf{g}_0$), and the final latent scale ($\boldsymbol{\ell}_0 = \boldsymbol{\ell}_{u,0} \odot \boldsymbol{\ell}_{v,0}$ by element-wise product), as shown in (Eqs. (7) and (8)):

$$\mathbf{U}_{\text{sign},0} = \text{sign}(\mathbf{U}'), \quad \mathbf{V}_{\text{sign},0} = \text{sign}(\mathbf{V}') \tag{6}$$

$$|\mathbf{U}'| \approx \mathbf{h}_0(\boldsymbol{\ell}_{u,0})^{\top}, \quad |\mathbf{V}'| \approx \mathbf{g}_0(\boldsymbol{\ell}_{v,0})^{\top} \tag{7}$$

$$\boldsymbol{\ell}_0 = \boldsymbol{\ell}_{u,0} \odot \boldsymbol{\ell}_{v,0} \tag{8}$$

The resulting initial effective weight $\widehat{\mathbf{W}}_{\text{pri},0} = \text{diag}(\mathbf{h}_0)\mathbf{U}_{\text{sign},0}\text{diag}(\boldsymbol{\ell}_0)\mathbf{V}_{\text{sign},0}^{\top}\text{diag}(\mathbf{g}_0)$ preserves key structural information. The learnable parameters are initialized as follows: the latent factors $\mathbf{U}$ and $\mathbf{V}$ are initialized with the factors $\mathbf{U}'$ and $\mathbf{V}'$ obtained from SVD respectively, and the scales $\mathbf{h}, \mathbf{g}, \boldsymbol{\ell}$ are initialized to $\mathbf{h}_0, \mathbf{g}_0, \boldsymbol{\ell}_0$.

### 3.3 Residual Compensation

To further improve fidelity without increasing the primary path's rank, we introduce **Residual Compensation**, a technique motivated by the theoretical insight that a two-stage error correction can be more effective than a single, larger approximation (Section A.1). Crucially, this method does not increase the model's total parameter budget; instead, it strategically reallocates the fixed bit budget of a single, higher-rank approximation into two lower-rank paths: a primary and a residual path. The term 'residual' primarily describes its role during initialization, where the auxiliary path is configured to model the error of the primary approximation. During QAT, both paths are optimized jointly to collectively represent the original weight. The advantage of this dual-path approach over a single path is empirically validated in our ablation studies (Section A.2.1).

To implement this, Residual Compensation employs a parallel auxiliary path that learns an approximation $\widehat{\mathbf{W}}_{\text{res}}$ of the residual error. This auxiliary path mirrors the structure of the LittleBit path:

$$\widehat{\mathbf{W}}_{\text{res}} = \text{diag}(\mathbf{h}_{\text{res}})\mathbf{U}_{\text{res,sign}}\text{diag}(\boldsymbol{\ell}_{\text{res}})\mathbf{V}_{\text{res,sign}}^{\top}\text{diag}(\mathbf{g}_{\text{res}}) \tag{9}$$

5

The parameters of this auxiliary path ($\mathbf{U}_{\mathrm{res,sign}}, \mathbf{V}_{\mathrm{res,sign}}, \mathbf{h}_{\mathrm{res}}, \mathbf{g}_{\mathrm{res}}, \boldsymbol{\ell}_{\mathrm{res}}$) are also learnable during QAT. They are initialized using the Dual-SVID strategy applied to the initial residual error $\mathbf{W}_{\mathrm{res,0}} = \mathbf{W} - \widehat{\mathbf{W}}_{\mathrm{pri,0}}$. Employing a separate path allows the model to specifically learn and compensate for the potentially distinct characteristics of this residual error. The final effective weight is the sum of both approximations:

$$\widehat{\mathbf{W}} = \widehat{\mathbf{W}}_{\mathrm{pri}} + \widehat{\mathbf{W}}_{\mathrm{res}} \tag{10}$$

This targeted error correction is often advantageous for maintaining performance at extreme compression levels, as demonstrated by the ablation studies presented in Section A.2.1. It is important to note that all learnable parameters from both the primary and residual paths are included when calculating the final BPW (see Section A.4 for calculation details).

Fig. 3 visualizes the initial state of weight components ($\widehat{\mathbf{W}}_{\mathrm{pri,0}}$, $\widehat{\mathbf{W}}_{\mathrm{res,0}}$, and their sum $\widehat{\mathbf{W}}_0$) for a Llama2-7B query weight, derived via Dual-SVID across effective bits from 0.1 to 1.0, compared against the original weight $\mathbf{W}$ (far right). Even at a low 0.1 BPW, the primary approximation $\widehat{\mathbf{W}}_{\mathrm{pri,0}}$ (top row, leftmost) captures the dominant low-rank structure of $\mathbf{W}$, albeit with considerable simplification. As effective bits increase, $\widehat{\mathbf{W}}_{\mathrm{pri,0}}$ progressively refines details and more effectively suppresses disruptive approximation noise, allowing underlying weight patterns to emerge with greater clarity. A key observation is that $\widehat{\mathbf{W}}_0$ (bottom row) at a modest 0.3 BPW—benefiting from the initialized residual component $\widehat{\mathbf{W}}_{\mathrm{res,0}}$ (middle row)—often presents a more faithful initial representation of $\mathbf{W}$ than $\widehat{\mathbf{W}}_{\mathrm{pri,0}}$ alone even at a higher 1.0 BPW. This underscores that Residual Compensation, even at the initialization stage, is vital for capturing complexities missed by the primary low-rank approximation, thereby providing a more faithful starting point for QAT.

## 4 Experiments

### 4.1 Settings

**Evaluation Setup** We evaluate LittleBit on diverse LLM families including Llama [48], Llama2 [49], Llama3 [50], OPT [51], Phi-4 [52], and QwQ [53], spanning sizes from 1.3B to 32B parameters. The primary performance metric is perplexity (PPL), evaluated on the WikiText-2 [54] validation dataset. Additional results on the C4 [55] and PTB [56] datasets are provided in Section A.3. We also assess zero-shot accuracy on common reasoning benchmarks: WinoGrande [61], OpenBookQA (OBQA) [59], HellaSwag [60], BoolQ [57], ARC-Easy (ARC-e), ARC-Challenge (ARC-c) [62], and PIQA [58].

**Training Details** The LittleBit model parameters, initialized using the Dual-SVID method (Section 3.2), are optimized using QAT with knowledge distillation (KD) [63, 36, 64]. The original pre-trained full-precision model serves as the teacher ($\mathcal{T}$) for the LittleBit student model ($\mathcal{S}$). The QAT objective combines the standard output Kullback-Leibler (KL) divergence loss $\mathcal{L}_{\mathrm{out}}$ and an intermediate layer Mean Squared Error (MSE) loss $\mathcal{L}_{\mathrm{inter}}$ to match hidden representations, weighted by an empirically determined $\lambda = 10$:

$$\mathcal{L}_{\mathrm{QAT}} = \mathcal{L}_{\mathrm{out}} + \lambda \mathcal{L}_{\mathrm{inter}}. \tag{11}$$

Following the setup in [24], the training data combine WikiText-2 and selected partitions from C4. We use 2048 token sequence length, 5 epochs, Adam optimizer ($\beta_1 = 0.9, \beta_2 = 0.999$) [65], and cosine LR decay with 2% warm-up (see details in Section A.7). For models employing Grouped-Query Attention (GQA) [66] (*e.g.*, Llama3, Phi-4, QwQ), the latent ranks $r$ for key(K) and value (V) projection layers were specifically adjusted (see Section A.2.3) to maintain performance under extreme quantization. For backpropagation, *SmoothSign* (forward: $\mathrm{sign}(x)$; backward gradient: $\tanh(kx)$, $k = 100$) is used, which generally performed better than Straight-Through Estimator (STE) [31] (see Section A.2.2). While the concept of a proxy gradient is not new, our contribution lies in the specific application and empirical validation of using the derivative of $\tanh(100x)$ as a stable gradient for QAT in the ultra-low-bit regime, which our ablation study confirms is effective (Section A.2.2).

**Baselines** We benchmark LittleBit across a range of effective bits from approximately 1.0 down to 0.1 bits per weight (BPW). In the challenging sub-1-bit regime, we compare against STBLLM [16],

Table 1: Perplexity (PPL) comparison on WikiText-2 across various LLMs and quantization methods. Lower PPL indicates better performance. BPW denotes effective bits per weight. Methods marked with † use quantization-aware training (QAT). STBLLM [16] utilizes N:M sparsity (ratio in parentheses). LittleBit demonstrates strong performance, particularly in the sub-1-bit regime.

| | Settings | | OPT | Llama | | Llama2 | | Llama3 | Phi-4 | QwQ |
|---|---|---|---|---|---|---|---|---|---|---|
| Method | Block Size | BPW | 1.3B | 7B | 13B | 7B | 13B | 8B | 14.7B | 32B |
| FullPrecision | - | 16 | 14.62 | 5.67 | 5.09 | 5.47 | 4.88 | 6.10 | 6.67 | 6.34 |
| OmniQuant | - | 2 | 28.82 | 9.75 | 7.83 | 11.20 | 8.25 | 349.27 | 12.09 | 10.19 |
| GPTQ | 128 | 2 | 119.98 | 39.96 | 15.08 | 52.22 | 23.63 | 1.5e3 | 24.96 | 67.16 |
| BiLLM | 128 | 1.1 | 74.07 | 41.95 | 13.95 | 29.00 | 23.71 | 54.29 | 16.95 | 15.4 |
| OneBit† | - | 1 | 20.27 | 8.21 | 7.37 | 8.36 | 7.41 | 13.09 | 9.92 | 9.86 |
| BinaryMoS† | - | 1 | 18.09 | 7.84 | 7.05 | 7.74 | 6.95 | 10.83 | 9.51 | 8.99 |
| LittleBit† | - | 1 | 20.33 | 9.03 | 8.17 | 9.08 | 8.18 | 16.30 | 11.28 | 12.08 |
| STBLLM | 128 | 0.80 (6:8) | 52.13 | 15.19 | 9.43 | 13.81 | 11.90 | 30.90 | 12.12 | 12.19 |
| STBLLM | 128 | 0.70 (5:8) | 73.42 | 19.52 | 11.47 | 19.17 | 14.00 | 59.83 | 14.57 | 13.78 |
| STBLLM | 128 | 0.55 (4:8) | 123.03 | 38.73 | 16.88 | 30.67 | 27.05 | 241.95 | 21.99 | 18.32 |
| STBLLM | 128 | 0.30 (2:8) | 2.3e3 | 1.6e3 | 592.06 | 1.8e3 | 893.82 | 1.7e5 | 761.05 | 512.01 |
| LittleBit† | - | 0.80 | 21.32 | 9.44 | 8.50 | 9.53 | 8.52 | 17.28 | 11.70 | 12.46 |
| LittleBit† | - | 0.70 | 21.99 | 9.66 | 8.71 | 9.85 | 8.76 | 18.01 | 12.05 | 13.22 |
| LittleBit† | - | 0.55 | 23.35 | 10.09 | 9.16 | 10.47 | 9.24 | 18.47 | 12.80 | 13.57 |
| LittleBit† | - | 0.30 | 28.30 | 11.50 | 10.33 | 12.00 | 10.48 | 20.34 | 14.71 | 16.48 |
| LittleBit† | - | 0.10 | 53.76 | 15.58 | 13.71 | 15.92 | 15.09 | 26.11 | 19.73 | 35.26 |

Table 2: Zero-shot accuracy (%) comparison on common sense reasoning benchmarks. Compares Full Precision (FP16) models against STBLLM [16] and LittleBit at 0.55 and 0.3 BPW.

| Models | Method | WinoGrande | OBQA | HellaSwag | BoolQ | ARC-e | ARC-c | PIQA | Average |
|---|---|---|---|---|---|---|---|---|---|
| Llama2-7B | FullPrecision | 67.16 | 40.80 | 72.95 | 71.37 | 69.44 | 40.95 | 78.12 | 62.97 |
| | STBLLM (0.55) | 52.80 | 33.00 | 36.94 | 62.17 | 37.33 | 25.34 | 62.46 | 44.29 |
| | STBLLM (0.3) | 50.43 | 31.80 | 26.20 | 37.83 | 25.88 | 25.09 | 53.26 | 35.78 |
| | LittleBit (0.55) | 51.62 | 34.00 | 44.57 | 61.38 | 44.15 | 26.96 | 68.12 | **47.26** |
| | LittleBit (0.3) | 51.30 | 34.80 | 37.61 | 61.80 | 39.98 | 25.09 | 65.83 | **45.20** |
| Llama2-13B | FullPrecision | 69.45 | 41.80 | 76.58 | 69.29 | 73.19 | 44.53 | 78.61 | 64.78 |
| | STBLLM (0.55) | 55.25 | 31.20 | 35.10 | 62.23 | 42.51 | 27.22 | 61.75 | 45.04 |
| | STBLLM (0.3) | 51.54 | 28.80 | 26.01 | 55.63 | 27.53 | 24.74 | 53.32 | 38.22 |
| | LittleBit (0.55) | 53.04 | 35.60 | 51.06 | 50.58 | 46.09 | 29.10 | 70.78 | **48.03** |
| | LittleBit (0.3) | 51.30 | 34.00 | 43.81 | 55.96 | 42.59 | 25.17 | 68.77 | **45.94** |
| Llama3-8B | FullPrecision | 72.92 | 45.00 | 79.18 | 81.25 | 80.21 | 52.98 | 79.54 | 70.15 |
| | STBLLM (0.55) | 52.17 | 25.80 | 30.61 | 57.16 | 30.35 | 23.72 | 57.56 | 39.62 |
| | STBLLM (0.3) | 49.57 | 26.60 | 26.36 | 51.62 | 26.05 | 24.40 | 51.74 | 36.62 |
| | LittleBit (0.55) | 50.12 | 30.20 | 36.78 | 57.55 | 46.80 | 22.95 | 66.27 | **44.38** |
| | LittleBit (0.3) | 51.93 | 28.20 | 33.91 | 57.98 | 43.48 | 24.32 | 64.64 | **43.49** |

a post-training quantization (PTQ) method employing N:M sparsity. For the ∼1.0 BPW setting, we compare against the 1-bit PTQ method BiLLM [7] and quantization-aware training (QAT) approaches (OneBit [13], BinaryMoS [24]). Results from standard low-bit methods like GPTQ [8] and OmniQuant [67] are included for broader context.

## 4.2 Main Results

**Superior Perplexity in Sub-1-bit Regime** Section 4.1 shows LittleBit's perplexity (PPL) results against baselines. LittleBit consistently performs robustly, especially in the sub-1-bit regime. Compared to STBLLM [16], LittleBit achieves markedly superior PPL at 0.8, 0.7, and 0.55 BPW. For instance, on Llama2-7B, LittleBit at 0.55 BPW achieves a PPL of 10.47, a substantial improvement over STBLLM's 30.67. Similar significant gains are observed across other models and BPWs; for Llama2-13B at 0.8 BPW, LittleBit records a PPL of 8.52 compared to STBLLM's 11.90.

**Extreme Low-BPW Stability** A key strength of LittleBit is its exceptional stability and performance in the extreme sub-0.5 BPW range, specifically at 0.3 BPW and 0.1 BPW. In this territory, methods like STBLLM exhibit severe performance degradation (*e.g.*, Llama2-7B PPL of 1.8e3 at 0.3 BPW). In contrast, LittleBit maintains strong PPL scores, such as 12.00 for Llama2-7B at 0.3 BPW and 15.92 at an unprecedented 0.1 BPW. When configured for an effective bits around 1.0 BPW by adjusting the latent rank $r$, LittleBit delivers perplexity comparable to specialized 1-bit QAT methods like OneBit [13] (*e.g.*, Llama2-7B: LittleBit 9.08 vs. OneBit 8.36) and PTQ methods

such as BiLLM [7]. While a method like BinaryMoS [24] may exhibit slightly lower PPL in some 1-bit scenarios (*e.g.*, Llama2-7B: BinaryMoS 7.74), this can be attributed to its use of more complex mechanisms like dynamic scaling, which represent a different set of architectural and computational trade-offs. Analyzing the trade-off by adjusting latent rank to tune effective bits from 1.0 down to 0.1 BPW reveals that while performance degradation is minimal down to 0.55 BPW, a *quantization cliff* appears between 0.3-0.1 BPW, positioning the 0.3-0.55 BPW range as an optimal *sweet spot* for balancing compression and accuracy.

**Model Generalization & Zero-Shot**    LittleBit's advantages are consistently observed across a diverse set of model architectures and sizes. This includes the Llama3 family, which is often considered challenging for quantization, and other contemporary models like Phi-4. For example, Llama3-8B quantized by LittleBit to 0.55 BPW yields a PPL of 18.47, significantly outperforming STBLLM's 241.95 (at 0.55 BPW with 4:8 sparsity).The zero-shot evaluation results presented in Section 4.1 further reinforce these perplexity-based findings. For instance, Llama2-7B compressed by LittleBit to 0.55 BPW achieves a mean accuracy of 47.26%, and even at an extreme 0.3 BPW, it maintains 45.20% mean accuracy. These results are notably competitive with, and in several cases surpass, those of STBLLM at similar or even higher effective bits (*e.g.*, STBLLM Llama2-7B at 0.55 BPW scores 44.29%). This suggests a superior preservation of the models' intrinsic reasoning capabilities even under extreme compression, highlighting LittleBit's efficacy in creating highly compact yet powerful LLMs.

## 4.3    Reasoning Performance of Extremely Compressed LLMs

The capability of LittleBit to achieve sub-1-bit compression while preserving significant model performance prompts an evaluation of their reasoning abilities, especially for large-scale models under extreme compression. We investigated the Phi-4 14B model, comparing its performance when compressed by LittleBit versus STBLLM across 7 zero-shot reasoning tasks. When compressed with LittleBit to 0.55 BPW, Phi-4 achieved an average accuracy of around 52.3%. While this performance is slightly lower than STBLLM at a similar compression level (0.55 BPW, 4:8 sparsity, achieving 54.2%), LittleBit demonstrates stronger graceful degradation. At a more aggressive 0.3 BPW, LittleBit-compressed Phi-4 maintained an average accuracy of 48.7%, significantly outperforming STBLLM at 0.3 BPW (2:8 sparsity), which scored 40.3%. Even at an extreme 0.1 BPW, Phi-4 compressed with LittleBit retained a notable average accuracy of 43.6%. These results highlight LittleBit's ability to maintain more robust reasoning capabilities compared to STBLLM as compression becomes increasingly extreme, particularly in the sub-0.5 BPW regime for the Phi-4 model.
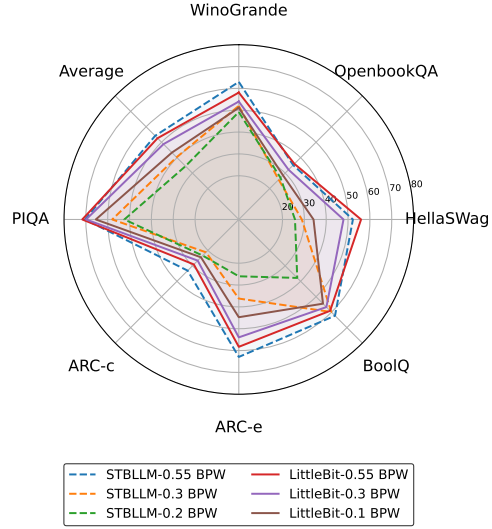


Figure 4: Zero-shot accuracy (%) on 7 common sense reasoning tasks for the Phi-4 14B model. Compares LittleBit-compressed Phi-4 at 0.55, 0.3, and 0.1 BPW against STBLLM.

## 5    Analysis

**Memory Footprint Reduction**    Designed for resource-constrained environments such as on-device deployment, LittleBit achieves substantial reductions in model memory footprint. As summarized in Section 5, applying LittleBit to quantize a Llama2-7B model (originally 13.49 GB in FP16) to an effective bit of 0.3 BPW reduces its required storage to 0.79 GB. This corresponds to a compression factor exceeding 17×. At the extreme setting of 0.1 BPW, the footprint is further reduced to 0.63 GB (over 21× compression). Comparable reductions are observed for larger models such as Llama2 70B, where the memory footprint decreases from 138.04 GB to under 2 GB at 0.1 BPW, achieving nearly 70× compression. These substantial memory reductions significantly expand the range of

Table 3: Memory footprint comparison (GB) for Llama2 [49] models under different quantization methods. LittleBit is evaluated at various BPWs, achieved by adjusting the latent rank $r$. Compression factors relative to FP16 are shown in parentheses.The gap between the effective BPW and the compression factor arises because quantization is applied only to the Transformer blocks, while components like the embedding layer and the lm_head remain in FP16.

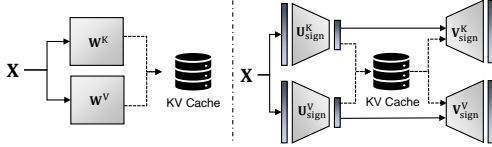| Model | FP16 | BiLLM [7] | OneBit [13] | LittleBit (Ours) | | | |
|---|---|---|---|---|---|---|---|
| BPW | 16 | 1.1 | 1 | 0.8 | 0.55 | 0.3 | 0.1 |
| Llama2-7B | 13.49 GB | 1.60 GB (8.43×) | 1.36 GB (9.92×) | 1.19 GB (11.34×) | 0.98 GB (13.77×) | 0.79 GB (17.08×) | 0.63 GB (21.41×) |
| Llama2-13B | 26.06 GB | 2.80 GB (9.31×) | 2.28 GB (11.43×) | 1.95 GB (13.36×) | 1.51 GB (17.26×) | 1.15 GB (22.66×) | 0.84 GB (31.02×) |
| Llama2-70B | 138.04 GB | 15.40 GB (8.96×) | 9.72 GB (14.20×) | 7.97 GB (17.31×) | 5.83 GB (23.68×) | 3.70 GB (37.31×) | 1.98 GB (69.72×) |



Figure 5: Conceptual view of KV Cache storage: the standard method (left) stores the full hidden dimension ($d_{\mathrm{model}}$), whereas LittleBit (right) caches a reduced latent dimension ($r$)

Table 4: Estimated KV cache memory reduction for Llama2-7B using LittleBit, with reduction factor $\sim d_{\mathrm{model}}/r$.

| BPW | KV Latent Rank ($r$) | KV Cache Reduction Factor |
|---|---|---|
| 0.80 | 1,624 | ~2.5× |
| 0.55 | 1,112 | ~3.7× |
| 0.30 | 600 | ~6.8× |
| 0.10 | 192 | ~21.3× |

devices capable of hosting large-scale language models, thereby enabling complex language tasks on hardware with limited VRAM or storage capacity.

**KV Cache Compression** LittleBit also inherently compresses the KV cache. Its factorization (Eq. (4)) applied to key/value projection matrices caches intermediate rank $r$ latent states (Fig. 5) instead of full $d_{\mathrm{model}}$-dimensional vectors, due to its forward computation (Eq. (5)). This intrinsically reduces KV cache memory requirements (especially for long contexts) by a factor proportional to $\sim d_{\mathrm{model}}/r$ (*e.g.*, up to 21.3× for Llama2-7B at 0.1 BPW, as shown inTable 4). This is similar to explicit methods (MLA [68], ASVD [42]) but achieved intrinsically via LittleBit's unified factorization in attention layers, concurrently reducing weight and activation memory.

**Latency Considerations** Beyond memory savings, inference speed is crucial. LittleBit's factorized architecture (Eq. (5)), unlike memory-focused or hardware-dependent methods, targets computational speedup via low-rank binary operations. To assess this potential, we benchmarked a custom 1-bit GEMV CUDA kernel designed to accelerate the primary computational path. As shown in Fig. 6 for a Llama2-70B MLP layer, the kernel achieves substantial speedups that increase as the bit-width decreases, reaching up to an **11.6×** speedup over a highly-optimized FP16 baseline at 0.1 BPW.
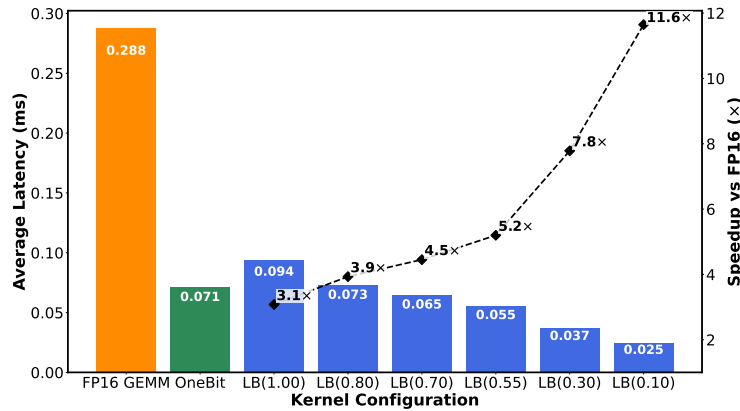


Figure 6: Kernel-level latency (bars, left axis) and speedup relative to FP16 (dashed line, right axis) on an NVIDIA A100 for a Llama2-70B MLP layer (8,192×28,672). The benchmark compares the FP16 GEMM baseline with OneBit and our LittleBit (LB) kernel at various BPW configurations from 1.0 to 0.1.

These results highlight LittleBit's potential for efficient deployment with significant latency reduction. The observed gap between theoretical computational gains and practical latency stems from the fact that inference at small batch sizes is often dominated by memory access rather than raw computation. Our factorized method involves multiple memory-bound steps, and our kernel, while effective, is not as exhaustively optimized as industry-standard libraries. A detailed analysis including theoretical FLOPs, practical latency, and end-to-end decoding throughput is provided in Section A.8.

## 6    Conclusion

In this work, we proposed LittleBit, a novel method advancing LLM compression to the extreme sub-0.5 BPW regime, with viable performance at 0.1 BPW. Gains were achieved via SVD-inspired latent matrix factorization, binarization of these factors, and a multi-scale (row, column, latent) compensation mechanism. We also introduced Dual-SVID initialization for stable quantization-aware training (QAT) and integrated Residual Compensation for error mitigation. Empirical results show LittleBit significantly outperforms prior sub-1-bit techniques and maintains robust fidelity across LLMs up to 32B parameters. LittleBit thus establishes a new state-of-the-art in the size–performance trade-off, significantly expanding LLM deployment in resource-constrained environments. Building on this foundation, future work can enhance practical deployment through hardware co-design for edge devices like neural processing units (NPUs). Further gains in model fidelity can be explored by integrating advanced techniques like non-uniform bit allocation, applied either across different model layers or even within a single layer by dynamically distributing the bit budget between the primary and residual paths.

## Acknowledgements

## References

[1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing System (NeurIPS)*, 30, 2017.

[2] Deepak Narayanan, Mohammad Shoeybi, Jared Casper, Patrick LeGresley, Mostofa Patwary, Vijay Anand Korthikanti, Dmitri Vainbrand, Prethvi Kashinkunti, Julie Bernauer, Bryan Catanzaro, et al. Efficient large-scale language model training on GPU clusters using Megatron-LM. *arXiv preprint arXiv:2104.04473*, 2021.

[3] Mikhail Isaev, Nic McDonald, and Richard Vuduc. Scaling infrastructure to support multi-trillion parameter LLM training. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC)*, 2023.

[4] Yue Zheng, Yuhao Chen, Bin Qian, Xiufang Shi, Yuanchao Shu, and Jiming Chen. A review on edge large language models: Design, execution, and applications. *arXiv preprint arXiv:2410.11845*, 2024.

[5] Xiurui Pan, Endian Li, Qiao Li, Shengwen Liang, Yizhou Shan, Ke Zhou, Yingwei Luo, Xiaolin Wang, and Jie Zhang. InstInfer: In-storage attention offloading for cost-effective long-context LLM inference. *arXiv preprint arXiv:2409.04992*, 2024.

[6] Libo Zhang, Zhaoning Zhang, Baizhou Xu, Songzhu Mei, and Dongsheng Li. Dovetail: A CPU/GPU heterogeneous speculative decoding for LLM inference. *arXiv preprint arXiv:2412.18934*, 2024.

[7] Wei Huang, Yangdong Liu, Haotong Qin, Ying Li, Shiming Zhang, Xianglong Liu, Michele Magno, and Xiaojuan Qi. BiLLM: pushing the limit of post-training quantization for LLMs. In *Proceedings of International Conference on Machine Learning (ICML)*, 2024.

[8] Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. GPTQ: Accurate post-training quantization for generative pre-trained transformers. In *Advances in Neural Information Processing Systems*, volume 36, pages 4424–4456, 2023.

[9] Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. AWQ: Activation-aware weight quantization for on-device LLM compression and acceleration. *Proceedings of Machine Learning and Systems (MLSys)*, 6:87–100, 2024.

[10] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2704–2713. IEEE, 2018.

[11] Jangho Kim, Yash Bhalgat, Jinwon Lee, Chirag Patel, and Nojun Kwak. QKD: Quantization-aware knowledge distillation. *arXiv preprint arXiv:1911.12491*, 2019.

[12] Hongyu Wang, Shuming Ma, Li Dong, Shaohan Huang, Huaijie Wang, Lingxiao Ma, Fan Yang, Ruiping Wang, Yi Wu, and Furu Wei. BitNet: Scaling 1-bit transformers for large language models. *arXiv preprint arXiv:2310.11453*, 2023.

[13] Yuzhuang Xu, Xu Han, Zonghan Yang, Shuo Wang, Qingfu Zhu, Zhiyuan Liu, Weidong Liu, and Wanxiang Che. OneBit: Towards extremely low-bit large language models. *Advances in Neural Information Processing System (NeurIPS)*, 37:66357–66382, 2024.

[14] Liqun Ma, Mingjie Sun, and Zhiqiang Shen. FBI-LLM: Scaling up fully binarized LLMs from scratch via autoregressive distillation. *arXiv preprint arXiv:2407.07093*, 2024.

[15] Yue Zheng, Yuhao Chen, Bin Qian, Xiufang Shi, Yuanchao Shu, and Jiming Chen. A review on edge large language models: Design, execution, and applications. *ACM Comput. Surv.*, 57(8), March 2025. ISSN 0360-0300. doi: 10.1145/3719664. URL https://doi.org/10.1145/3719664.

[16] Peijie Dong, Lujun Li, Yuedong Zhong, Dayou Du, Ruibo Fan, Yuhan Chen, Zhenheng Tang, Qiang Wang, Wei Xue, Yike Guo, et al. STBLLM: Breaking the 1-bit barrier with structured binary LLMs. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2025.

[17] Ajay Jaiswal, Lu Yin, Zhenyu Zhang, Shiwei Liu, Jiawei Zhao, Yuandong Tian, and Zhangyang Wang. From galore to welore: How low-rank weights non-uniformly emerge from low-rank gradients. *arXiv preprint arXiv:2407.11239*, 2024.

[18] Qinsi Wang, Jinghan Ke, Masayoshi Tomizuka, Yiran Chen, Kurt Keutzer, and Chenfeng Xu. Dobi-SVD: Differentiable SVD for LLM compression and some new perspectives. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2025.

[19] Gene H. Golub and Christian Reinsch. Singular value decomposition and least squares solutions. *Numerische Mathematik*, 14(5):403–420, 1970. doi: 10.1007/BF02163027.

[20] Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2016. URL https://arxiv.org/abs/1510.00149.

[21] Xin Wang, Samiul Alam, Zhongwei Wan, Hui Shen, and Mi Zhang. SVD-LLM V2: Optimizing singular value truncation for large language model compression. *Annual Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics (NAACL)*, 2025.

[22] Mingjie Sun, Zhuang Liu, Anna Bair, and J. Zico Kolter. A simple and effective pruning approach for large language models. In *Proceedings of the Twelfth International Conference on Learning Representations (ICLR)*, 2024.

[23] Ruihao Gong, Yifu Ding, Zining Wang, Chengtao Lv, Xingyu Zheng, Jinyang Du, Haotong Qin, Jinyang Guo, Michele Magno, and Xianglong Liu. A survey of low-bit large language models: Basics, systems, and algorithms. *arXiv preprint arXiv:2409.16694*, 2024.

[24] Dongwon Jo, Taesu Kim, Yulhwa Kim, and Jae-Joon Kim. Mixture of Scales: Memory-efficient token-adaptive binarization for large language models. *Advances in Neural Information Processing System (NeurIPS)*, 37:137474–137494, 2024.

[25] Zhiteng Li, Xianglong Yan, Tianao Zhang, Haotong Qin, Dong Xie, Jiang Tian, Linghe Kong, Yulun Zhang, Xiaokang Yang, et al. ARB-LLM: Alternating refined binarizations for large language models. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2025.

[26] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2016.

[27] Cheng Fu, Shilin Zhu, Hao Su, Ching-En Lee, and Jishen Zhao. Towards fast and energy-efficient binarized neural network inference on fpga. *arXiv preprint arXiv:1810.02068*, 2018.

[28] Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. BinaryConnect: Training deep neural networks with binary weights during propagations. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2015.

[29] Haotong Qin, Ruihao Gong, Xianglong Liu, Xiao Bai, Jingkuan Song, and Nicu Sebe. Binary neural networks: A survey. *Pattern Recognition*, 105:107281, 2020.

[30] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. XNOR-Net: Imagenet classification using binary convolutional neural networks. In *European conference on computer vision*, pages 525–542. Springer, 2016.

[31] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.

[32] Matt Schoenbauer, Daniele Moro, Lukasz Lew, and Andrew Howard. Custom gradient estimators are straight-through estimators in disguise. *arXiv preprint arXiv:2405.05171*, 2024.

[33] Yuzhang Shang, Zhihang Yuan, Qiang Wu, and Zhen Dong. PB-LLM: Partially binarized large language models. *arXiv preprint arXiv:2310.00034*, 2023.

[34] Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. A survey of quantization methods for efficient neural network inference. *arXiv preprint arXiv:2103.13630*, 2021. URL `https://arxiv.org/abs/2103.13630`.

[35] Misha Denil, Babak Shakibi, Laurent Dinh, Marc'Aurelio Ranzato, and Nando de Freitas. Predicting parameters in deep learning. In *Advances in Neural Information Processing Systems*, volume 26, pages 2148–2156, 2013. URL `https://papers.nips.cc/paper/5025-predicting-parameters-in-deep-learning`.

[36] Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang Shi, Raghuraman Krishnamoorthi, and Vikas Chandra. LLM-QAT: Data-free quantization aware training for large language models. In *Findings of the Association for Computational Linguistics (ACL)*, pages 467–484, 2024.

[37] Xin Men, Mingyu Xu, Qingyu Zhang, Bingning Wang, Hongyu Lin, Yaojie Lu, Xianpei Han, and Weipeng Chen. ShortGPT: Layers in large language models are more redundant than you expect. *arXiv preprint arXiv:2403.03853*, 2024. URL `https://arxiv.org/abs/2403.03853`.

[38] Shwai He, Guoheng Sun, Zheyu Shen, and Ang Li. What matters in transformers? not all attention is needed. *arXiv preprint arXiv:2406.15786*, 2024.

[39] Rajarshi Saha, Naomi Sagan, Varun Srivastava, Andrea Goldsmith, and Mert Pilanci. Compressing large language models using low rank and low precision decomposition. In *Advances in Neural Information Processing Systems (NeurIPS) 2024 Poster Session*, 2024. URL `https://openreview.net/forum?id=lkx30pcqSZ`.

[40] Haokun Zhang, Xian Li, Yuxuan Li, Yiming Liu, Yuxiang Wang, and Yisen Liu. Adaptive budget allocation for parameter-efficient fine-tuning. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2023. URL `https://arxiv.org/abs/2303.10512`.

[41] Haokun Zhang, Xian Li, Yuxuan Li, Yiming Liu, Yuxiang Wang, and Yisen Liu. PiSSA: Principal singular values and singular vectors adaptation of pre-trained models. *Advances in Neural Information Processing Systems*, 2024. URL `https://arxiv.org/abs/2404.02948`.

[42] Zhihang Yuan, Yuzhang Shang, Yue Song, Qiang Wu, Yan Yan, and Guangyu Sun. ASVD: Activation-aware singular value decomposition for compressing large language models. *arXiv preprint arXiv:2312.05821*, 2023.

[43] Cheng Zhang, Jianyi Cheng, George Anthony Constantinides, and Yiren Zhao. LQER: Low-rank quantization error reconstruction for LLMs. In *Proceedings of the 41st International Conference on Machine Learning*, pages 58763–58779. PMLR, 2024. URL `https://proceedings.mlr.press/v235/zhang24j.html`.

[44] Wenjing Ke, Zhe Li, Dong Li, Lu Tian, and Emad Barsoum. DL-QAT: Weight-decomposed low-rank quantization-aware training for large language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing (Industry Track)*, pages 85–94. Association for Computational Linguistics, 2024. URL `https://aclanthology.org/2024.emnlp-industry.10/`.

[45] Xin Wang, Yu Zheng, Zhongwei Wan, and Mi Zhang. SVD-LLM: Truncation-aware singular value decomposition for large language model compression. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2025.

[46] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.

[47] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, 2010.

[48] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. LLaMA: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[49] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

[50] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The Llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

[51] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. OPT: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.

[52] Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, et al. Phi-4 technical report. *arXiv preprint arXiv:2412.08905*, 2024.

[53] Qwen Team. QwQ-32b: Embracing the power of reinforcement learning, March 2025. URL `https://qwenlm.github.io/blog/qwq-32b/`.

[54] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. In *Proceedings of the Fifth International Conference on Learning Representations (ICLR)*, 2017.

[55] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.

[56] Mitch Marcus, Grace Kim, Mary Ann Marcinkiewicz, Robert MacIntyre, Ann Bies, Mark Ferguson, Karen Katz, and Britta Schasberger. The Penn Treebank: Annotating predicate argument structure. In *Human Language Technology: Proceedings of a Workshop held at Plainsboro, New Jersey, March 8-11, 1994*, 1994.

[57] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 2924–2936, 2019.

[58] Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. PIQA: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence (AAAI)*, pages 7432–7439, 2020.

[59] Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *EMNLP*, 2018.

[60] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 4791–4800, 2019.

[61] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. WinoGrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.

[62] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? Try ARC, the AI2 Reasoning Challenge. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1168–1179. Association for Computational Linguistics, 2018. URL `https://aclanthology.org/D18-1149/`.

[63] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.

[64] Wei Zhang, Lu Hou, Yichun Yin, Lifeng Shang, Xiao Chen, Xin Jiang, and Qun Liu. TernaryBERT: Distillation-aware ultra-low bit BERT. *arXiv preprint arXiv:2009.12812*, 2020.

[65] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of the Third International Conference on Learning Representations (ICLR)*, 2015.

[66] Joshua Ainslie, James Lee-Thorp, Michiel De Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit Sanghai. GQA: Training generalized multi-query transformer models from multi-head checkpoints. *arXiv preprint arXiv:2305.13245*, 2023.

[67] Wenqi Shao, Mengzhao Chen, Zhaoyang Zhang, Peng Xu, Lirui Zhao, Zhiqian Li, Kaipeng Zhang, Peng Gao, Yu Qiao, and Ping Luo. OmniQuant: Omnidirectionally calibrated quantization for large language models. In *Proceedings of the Twelfth International Conference on Learning Representations (ICLR)*, 2024.

[68] DeepSeek-AI. DeepSeek-V2: A strong mixture-of-experts language model with expert attention. *arXiv preprint arXiv:2405.04434*, 2024.

[69] Noam Shazeer. Fast transformer decoding: One write-head is all you need. *arXiv preprint arXiv:1911.02150*, 2019.

[70] Casper Hansen. AutoAWQ kernels. `https://github.com/casper-hansen/AutoAWQ_kernels`, 2024. Accessed: 2025-04-01.

# A  Appendix

## A.1  Mathematical Analysis and Proofs

This section provides mathematical details for the LittleBit computation and explores theoretical aspects related to low-rank quantization errors, offering motivation for certain design choices.

**Proposition 1** (LittleBit Forward Pass Computation (Proof Detail)). *As stated in Section 3.1 (Proposition 1), for input $\mathbf{X} \in \mathbb{R}^{\text{seq} \times d_{\text{in}}}$, the primary quantized weight matrix is $\widehat{\mathbf{W}}_{\text{pri}} = \text{diag}(\mathbf{h})\mathbf{U}_{\text{sign}}\text{diag}(\boldsymbol{\ell})\mathbf{V}_{\text{sign}}^{\top}\text{diag}(\mathbf{g})$. The forward pass $\mathbf{Y} = \mathbf{X}\widehat{\mathbf{W}}_{\text{pri}}^{\top}$ can be computed as shown in Eq. (12).*

$$\mathbf{Y} = ((((\mathbf{X} \odot \mathbf{g})\mathbf{V}_{\text{sign}}) \odot \boldsymbol{\ell})\mathbf{U}_{\text{sign}}^{\top}) \odot \mathbf{h}. \tag{12}$$

*Terms are defined in Section 3.1.*

*Proof.* Starting with $\mathbf{Y} = \mathbf{X}\widehat{\mathbf{W}}_{\text{pri}}^{\top}$:

$$\begin{aligned}
\mathbf{Y} &= \mathbf{X} \left( \text{diag}(\mathbf{h})\mathbf{U}_{\text{sign}}\text{diag}(\boldsymbol{\ell})\mathbf{V}_{\text{sign}}^{\top}\text{diag}(\mathbf{g}) \right)^{\top} \\
&= \mathbf{X}\,\text{diag}(\mathbf{g})\mathbf{V}_{\text{sign}}\text{diag}(\boldsymbol{\ell})\mathbf{U}_{\text{sign}}^{\top}\text{diag}(\mathbf{h}) \\
&= ((((\mathbf{X} \odot \mathbf{g})\mathbf{V}_{\text{sign}}) \odot \boldsymbol{\ell})\mathbf{U}_{\text{sign}}^{\top}) \odot \mathbf{h}.
\end{aligned} \tag{13}$$

This yields Eq. (12), where $\odot$ implies element-wise multiplication with broadcasting. $\square$

**Claim 1** (Quantization Error vs. Factor Rank). *Consider a simplified model where a rank-$r$ matrix $\mathbf{W}^{(r)} = \mathbf{U}^{(r)}(\mathbf{V}^{(r)})^{\top}$ is approximated. Let $\mathbf{U}_{\text{sign}}^{(r)} = \text{sign}(\mathbf{U}^{(r)})$, $\mathbf{V}_{\text{sign}}^{(r)} = \text{sign}(\mathbf{V}^{(r)})$. If magnitudes are crudely approximated using fixed rank-1 SVDs of $|\mathbf{U}^{(r)}|$ and $|\mathbf{V}^{(r)}|$ to form scales $\mathbf{s}_U^{(r)}, \mathbf{s}_V^{(r)}$:*

$$\widehat{\mathbf{W}}_r := \left( \mathbf{U}_{\text{sign}}^{(r)} \odot (\mathbf{s}_U^{(r)}\mathbf{1}_r^{\top}) \right) \left( \mathbf{V}_{\text{sign}}^{(r)} \odot (\mathbf{s}_V^{(r)}\mathbf{1}_r^{\top}) \right)^{\top}. \tag{14}$$

*The error $\mathcal{E}(r) := \left\| \mathbf{W}^{(r)} - \widehat{\mathbf{W}}_r \right\|_F$ might be non-decreasing with $r$.*

*Proof Sketch / Heuristic Argument.* As rank $r$ (and thus complexity of $\mathbf{W}^{(r)}, \mathbf{U}^{(r)}, \mathbf{V}^{(r)}$) increases, the fixed-structure rank-1 magnitude approximation (via $\mathbf{s}_U^{(r)}, \mathbf{s}_V^{(r)}$) may become a relatively poorer fit for the increasingly complex actual magnitudes of $\mathbf{U}^{(r)}, \mathbf{V}^{(r)}$. This degradation in scaling accuracy could lead to $\mathcal{E}(r)$ not decreasing, or even increasing, motivating more sophisticated scaling as in LittleBit. $\square$

**Claim 2** (Quantization Bias vs. SVD Component Structure). *Let $\mathbf{W}$'s rank-$r$ SVD approximation be $\widetilde{\mathbf{W}}_r = \widetilde{\mathbf{W}}_{r_1} + \widetilde{\mathbf{R}}_{r_2}$ (primary $\widetilde{\mathbf{W}}_{r_1}$, residual $\widetilde{\mathbf{R}}_{r_2}$). Consider a quantization operator $\mathcal{Q}$. The relative quantization error $\|\mathcal{Q}(\mathbf{A}) - \mathbf{A}\|_F / \|\mathbf{A}\|_F$ might be larger for matrices $\mathbf{A}$ with "flatter" singular value spectra or less energy concentration (e.g., $\widetilde{\mathbf{R}}_{r_2}$ vs. $\widetilde{\mathbf{W}}_{r_1}$). This suggests that $\mathcal{Q}$ applied to $\widetilde{\mathbf{R}}_{r_2}$ might be less effective (i.e., higher relative quantization error) than when applied to $\widetilde{\mathbf{W}}_{r_1}$.*

*Proof Sketch.* $\widetilde{\mathbf{W}}_{r_1}$ has dominant singular values, while $\widetilde{\mathbf{R}}_{r_2}$ has smaller ones, often resulting in lower energy concentration for $\widetilde{\mathbf{R}}_{r_2}$. If $\mathcal{Q}$'s relative error is sensitive to this (*e.g.*, higher for lower concentration), quantizing $\widetilde{\mathbf{W}}_{r_1}$ and $\widetilde{\mathbf{R}}_{r_2}$ separately may be better than quantizing $\widetilde{\mathbf{W}}_r$ jointly, as it allows adapting $\mathcal{Q}$ to their differing structures. $\square$

**Proposition 2** (Potential Advantage of Two-Stage Quantization of SVD Components). *Using the setup of Claim 2, let $\widehat{\mathbf{W}}_r^{\text{single}} = \mathcal{Q}(\widetilde{\mathbf{W}}_r)$ and $\widehat{\mathbf{W}}^{\text{two-stage}} = \mathcal{Q}(\widetilde{\mathbf{W}}_{r_1}) + \mathcal{Q}(\widetilde{\mathbf{R}}_{r_2})$. The two-stage error $\mathcal{E}_{\text{two-stage}} = \left\| \mathbf{W} - \widehat{\mathbf{W}}^{\text{two-stage}} \right\|_F$ can be less than the single-stage error $\mathcal{E}_{\text{single-stage}} = \left\| \mathbf{W} - \widehat{\mathbf{W}}_r^{\text{single}} \right\|_F$. This occurs if the quantization error of the sum is greater than the sum of (vectorial) quantization errors of parts:*

$$\|(\mathcal{Q}(\widetilde{\mathbf{W}}_{r_1}) - \widetilde{\mathbf{W}}_{r_1}) + (\mathcal{Q}(\widetilde{\mathbf{R}}_{r_2}) - \widetilde{\mathbf{R}}_{r_2})\|_F < \|\mathcal{Q}(\widetilde{\mathbf{W}}_r) - \widetilde{\mathbf{W}}_r\|_F. \tag{15}$$

15

*Proof Outline and Discussion.* Let $\mathbf{E}_{\text{trunc}} = \mathbf{W} - \widetilde{\mathbf{W}}_r$. Errors are $\Delta_1 = \mathcal{Q}(\widetilde{\mathbf{W}}_{r_1}) - \widetilde{\mathbf{W}}_{r_1}$, $\Delta_2 = \mathcal{Q}(\widetilde{\mathbf{R}}_{r_2}) - \widetilde{\mathbf{R}}_{r_2}$, $\Delta_{\text{sum}} = \mathcal{Q}(\widetilde{\mathbf{W}}_r) - \widetilde{\mathbf{W}}_r$. Then $\mathcal{E}_{\text{two}-\text{stage}} = \|\mathbf{E}_{\text{trunc}} - (\Delta_1 + \Delta_2)\|_F$ and $\mathcal{E}_{\text{single}-\text{stage}} = \|\mathbf{E}_{\text{trunc}} - \Delta_{\text{sum}}\|_F$. The condition Eq. (15) makes $\mathcal{E}_{\text{two}-\text{stage}} < \mathcal{E}_{\text{single}-\text{stage}}$ likely. Non-linear $\mathcal{Q}$ (*e.g.*, with adaptive scaling) can satisfy this by better handling the distinct characteristics of $\widetilde{\mathbf{W}}_{r_1}$ and $\widetilde{\mathbf{R}}_{r_2}$. This motivates LittleBit's separate residual handling, though its residual is learned differently. $\qquad\square$

## A.2 Ablation Study

### A.2.1 Residual Compensation

To provide further intuition behind the benefits of Residual Compensation, we visualize the output activations of selected Transformer layers under different quantization schemes. Specifically, we extract the final outputs from Transformer layers 0, 5, 10, and 15 of an OPT-1.3B model, evaluated on two randomly sampled input sequences. For each layer and input, we display a set of three activation maps corresponding to the full-precision baseline, as well as two quantized variants—with and without residual compensation. These groupings allow us to directly observe the structural differences induced by the presence or absence of the residual compensation.

As shown in Fig. 7, residual compensation preserves a high-fidelity resemblance to the full-precision baseline across all inspected layers. In the no-residual setting, the output feature distributions exhibit noticeable deformation relative to the full-precision baseline. In contrast, applying residual compensation consistently restores structural similarity to the original distribution, preserving both directional patterns and magnitude diversity.

 Table 5 presents the PPL results for the OPT-1.3B model, comparing configurations with and without residual compensation. As indicated in the table, at the extremely low bit of 0.1 BPW, the version with residual compensation (PPL 60.011) performed worse than the version without it (PPL 48.512). This suggests that for relatively small models (1.3B), when pushing effective bits to such extreme lows, the additional complexity introduced by the residual compensation mechanism might outweigh its benefits and instead hinder performance. However, for the same OPT-1.3B model at other effective bits evaluated (ranging from 0.3 BPW to 1.0 BPW), the inclusion of residual compensation consistently resulted in better (lower) PPL. Moreover, experiments conducted with other, larger models generally showed that residual compensation tended to provide superior performance. In light of these overall advantages observed across various settings and particularly for larger models, residual compensation was adopted as a standard component in all experiments presented in this study.

Table 5: Perplexity (WikiText-2) comparison between Residual and Non-Residual for OPT-1.3B at various BPWs. Learning rate was 8e-5. Lower PPL is better.

| BPW | Residual PPL | Non-Residual PPL | Difference (Residual - Non-Residual) |
|-----|-----|-----|-----|
| 1.00 | 20.329 | 21.691 | -1.362 |
| 0.80 | 21.338 | 22.688 | -1.350 |
| 0.70 | 22.001 | 23.313 | -1.312 |
| 0.55 | 23.457 | 24.788 | -1.331 |
| 0.30 | 28.984 | 29.724 | -0.740 |
| 0.10 | 60.011 | 48.512 | 11.499 |

### A.2.2 SmoothSign versus Straight-Through Estimator

For backpropagating through the non-differentiable $\text{sign}(x)$ function, we compared the Straight-Through Estimator (STE) [31] with our proposed *SmoothSign* technique. SmoothSign employs the $\text{sign}(x)$ function for the forward pass. For the backward pass, it utilizes a smooth proxy gradient, specifically the derivative of $\tanh(kx)$ where $k = 100$ (illustrated in Fig. 8).

Table 6 (OPT-1.3B results) shows SmoothSign yielding better PPL as effective bits decrease, with a notable advantage at 0.1 BPW. Due to its superior stability and performance in the ultra-low bit regime, SmoothSign was adopted for all QAT experiments (consistent with Section 4.1).

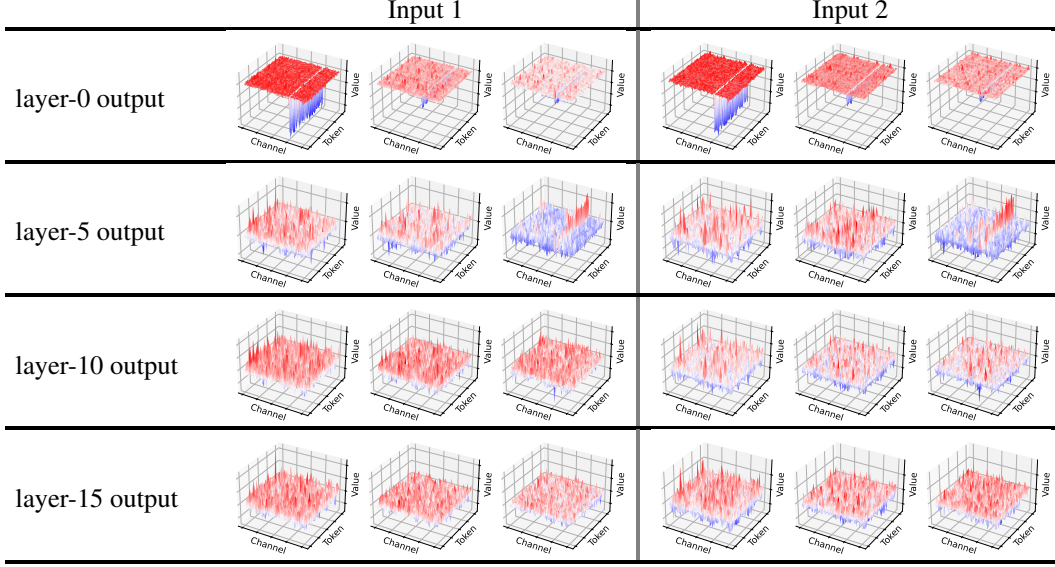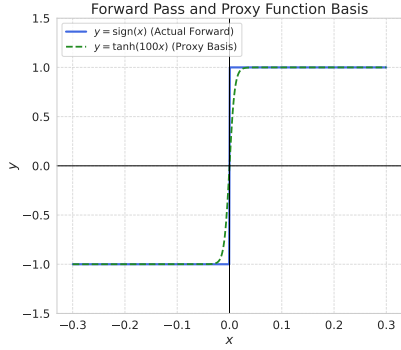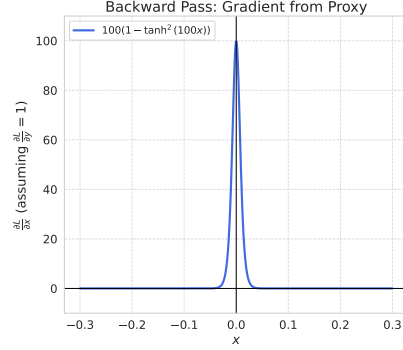|  | Input 1 | Input 2 |
|---|---|---|
| layer-0 output | | |
| layer-5 output | | |
| layer-10 output | | |
| layer-15 output | | |

Figure 7: Visualization of Transformer layer outputs for layers 0, 5, 10, and 15 under the LittleBit quantization scheme on OPT-1.3B. For each of two randomly sampled input sequences (annotated at the top), we show horizontally grouped triplets of activation maps per layer: full-precision baseline (left), quantized with residual compensation (center), and quantized without residual compensation (right).



(a) Forward pass: $\text{sign}(x)$ and $\tanh(100x)$

(b) Backward pass: gradient $\frac{d}{dx}\tanh(100x)$

Figure 8: The SmoothSign activation technique. (a) Functions relevant to the forward pass: Smooth-Sign uses the $\text{sign}(x)$ function (actual forward pass). The plot also shows $\tanh(100x)$, which serves as the basis for deriving the smooth proxy gradient shown in (b). (b) The smooth proxy gradient, $\frac{d}{dx}\tanh(100x)$, utilized for the backward pass.

### A.2.3 Grouped Query Attention

Modern LLMs often use Grouped Query Attention (GQA) [66] or Multi-Query Attention (MQA) [69], where key (K) and value (V) projection layers have smaller output dimensions than query (Q) projections. Applying LittleBit at ultra-low bits (*e.g.*, 0.1 BPW) to these GQA/MQA models can result in an extremely small latent dimension $r$ for K/V layers (*e.g.*, $r \approx 20$ for Llama3-8B K/V layers at 0.1 BPW overall), potentially creating an information bottleneck and hindering QAT or performance.

To investigate, we performed an ablation on Llama3-8B ($\sim$0.1 BPW target), increasing the latent rank $r$ for K/V layers by 2$\times$, 4$\times$, and 8x over the standard calculation (Section A.4), while other layers remained unchanged. Table 7 shows that a 4x factor for K/V rank yielded the best PPL on WikiText-2

Table 6: Perplexity (WikiText-2) comparison between STE and SmoothSign for OPT-1.3B at various effective bits (BPW). Learning rate was 8e-5. Lower PPL is better.

| BPW | STE PPL | SmoothSign PPL | Difference (SmoothSign - STE) |
|------|---------|----------------|-------------------------------|
| 1.00 | 20.322 | 20.329 | +0.007 |
| 0.80 | 21.344 | 21.338 | -0.006 |
| 0.70 | 22.100 | 22.001 | -0.099 |
| 0.55 | 23.528 | 23.457 | -0.071 |
| 0.30 | 29.058 | 28.984 | -0.074 |
| 0.10 | 60.401 | 60.011 | -0.390 |

Table 7: Ablation study on adjusting the latent dimension factor ($r$) for key (K) and value (V) projection layers in Llama3-8B with GQA, targeting an overall effective bit of 0.1 BPW. Performance measured by perplexity (PPL, lower is better) on WikiText-2 and C4 validation sets, and average zero-shot accuracy (%) on common sense reasoning tasks (higher is better). The 'KV Factor' indicates the multiplier applied to the K/V layers' latent rank $r$ compared to the standard calculation for 0.1 BPW.

| KV Factor | Approx. BPW | WikiText-2 PPL | C4 PPL | Avg. Zero-shot Acc (%) |
|-----------|-------------|----------------|--------|------------------------|
| $1\times$ (Baseline) | $\sim$0.098 | 25.80 | 34.04 | 44.85 |
| $2\times$ | $\sim$0.101 | 25.22 | 33.20 | **44.92** |
| $4\times$ | $\sim$0.107 | **24.93** | **32.69** | 44.70 |
| $8\times$ | $\sim$0.119 | 25.31 | 33.21 | 44.87 |

and C4, with minimal impact on overall BPW (*e.g.*, $\sim$0.098 BPW to $\sim$0.107 BPW for $4\times$). Average zero-shot accuracy also remained competitive. An 8x factor showed diminishing returns.

Given the significant PPL improvement with the $4\times$ K/V rank factor (nearly 1 point on WikiText-2 vs. $1\times$) at a negligible BPW cost (<10% relative increase in effective bits for transformer layers, <4% of total parameters), this strategy was adopted for GQA/MQA models (Llama3-8B, Phi-4, QwQ-32B) in our main results. For these models, the latent rank $r$ for K and V projection layers was calculated based on a $4\times$ multiplier compared to the standard calculation for the target BPW, while all other layers used the standard rank calculation. This approach helps mitigate potential information bottlenecks in GQA/MQA layers under extreme quantization without significantly altering the target compression level.

## A.3  Perplexity on C4 and PTB

To further assess LittleBit's generalization, we evaluated its performance on additional benchmarks: the diverse C4 web text corpus and the distinct PTB news corpus (Table 8). Focusing on Llama2-7B/13B against STBLLM, the results on these datasets confirm the trends observed on WikiText-2. LittleBit consistently achieves lower perplexity than STBLLM at comparable sub-1-bit regimes. This performance advantage becomes more significant below 0.55 BPW, where STBLLM's performance degrades considerably, particularly on PTB.

Notably, LittleBit maintains robust performance and stability even at the extreme 0.1 BPW level across all tested datasets. This consistent ability to maintain strong language modeling capabilities under severe compression highlights the generalizability of our approach, which integrates latent factorization, multi-scale compensation, and QAT. It reinforces LittleBit's potential for effectively deploying capable LLMs in resource-limited environments.

## A.4  Average Bits Per Weight Calculation

This section details the calculation of the average effective bits per weight (which we denote as $b$ in the following equations) for a linear layer implemented using the LittleBit architecture, including the parameters from the primary ($\widehat{\mathbf{W}}_{\text{pri}}$) and residual compensation ($\widehat{\mathbf{W}}_{\text{res}}$) pathways, as described in Sections 3.1 and 3.3. The calculation demonstrates how to determine the latent rank $r$ required to achieve a target value for $b$ within the range of approximately 0.1 to 1.0.

Table 8: Perplexity (PPL) comparison on C4 and PTB validation sets for Llama2 models. Lower PPL indicates better performance. STBLLM [16] utilizes N:M sparsity (ratio in parentheses).

| Model | Method | BPW | C4 | PTB |
|-------|--------|-----|-----|-----|
| Llama2-7B | FullPrecision | 16 | 6.97 | 37.91 |
| | STBLLM (6:8) | 0.80 | 15.42 | 2.4e3 |
| | STBLLM (4:8) | 0.55 | 30.99 | 527.56 |
| | STBLLM (2:8) | 0.30 | 808.98 | 2.3e4 |
| | LittleBit | 0.80 | 11.77 | 40.75 |
| | LittleBit | 0.55 | 12.94 | 50.53 |
| | LittleBit | 0.30 | 14.78 | 59.17 |
| | LittleBit | 0.10 | 19.73 | 83.09 |
| Llama2-13B | FullPrecision | 16 | 6.47 | 50.94 |
| | STBLLM (6:8) | 0.80 | 12.96 | 165.51 |
| | STBLLM (4:8) | 0.55 | 27.38 | 424.82 |
| | STBLLM (2:8) | 0.30 | 442.49 | 1.7e3 |
| | LittleBit | 0.80 | 10.60 | 44.87 |
| | LittleBit | 0.55 | 11.53 | 53.52 |
| | LittleBit | 0.30 | 13.07 | 51.66 |
| | LittleBit | 0.10 | 18.88 | 87.59 |

A LittleBit linear layer replaces the original weight matrix $\mathbf{W} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ with two parallel structures, each comprising:

- Two binary sign matrices: $\mathbf{U}_{\text{sign}} \in \{\pm 1\}^{d_{\text{out}} \times r}$ and $\mathbf{V}_{\text{sign}} \in \{\pm 1\}^{d_{\text{in}} \times r}$. These require 1 bit per element.

- Three FP16 scaling vectors: $\mathbf{h} \in \mathbb{R}^{d_{\text{out}}}$, $\mathbf{g} \in \mathbb{R}^{d_{\text{in}}}$, and $\boldsymbol{\ell} \in \mathbb{R}^r$. These require 16 bits per element.

Since both the primary and residual paths have this structure, the total number of bits required to store the parameters for a LittleBit layer is:

$$\begin{aligned}
\text{Total Bits} &= 2 \times (\text{Bits for } \mathbf{U}_{\text{sign}} + \text{Bits for } \mathbf{V}_{\text{sign}} + \text{Bits for } \mathbf{h} + \text{Bits for } \mathbf{g} + \text{Bits for } \boldsymbol{\ell}) \\
&= 2 \times ((d_{\text{out}} \times r \times 1) + (d_{\text{in}} \times r \times 1) + (d_{\text{out}} \times 16) + (d_{\text{in}} \times 16) + (r \times 16)) \\
&= 2r(d_{\text{out}} + d_{\text{in}}) + 32(d_{\text{out}} + d_{\text{in}}) + 32r
\end{aligned}$$

Note: If Residual Compensation is not used, the initial factor of 2 should be removed from the calculation.

The average bits per weight, denoted as $b$, is calculated by dividing the total bits by the number of parameters in the original FP16 weight matrix ($d_{\text{out}} \times d_{\text{in}}$):

$$b = \frac{2r(d_{\text{out}} + d_{\text{in}}) + 32(d_{\text{out}} + d_{\text{in}}) + 32r}{d_{\text{out}} \times d_{\text{in}}} \tag{16}$$

To achieve a target value for $b$, we rearrange Eq. (16) to solve for the required latent rank $r$:

$$r = \frac{(b \times d_{\text{out}} \times d_{\text{in}}) - 32(d_{\text{out}} + d_{\text{in}})}{2(d_{\text{out}} + d_{\text{in}}) + 32} \tag{17}$$

Since $r$ must be an integer, we typically round the calculated value to the nearest suitable integer. This chosen integer $r$ then determines the actual value of $b$ achieved, which will be very close to the target value.

- **Example 1: Linear Layer** ($4{,}096 \times 4{,}096$) Let $d_{\text{out}} = 4{,}096$ and $d_{\text{in}} = 4{,}096$. The original number of parameters is $4{,}096 \times 4{,}096 = 16{,}777{,}216$. $d_{\text{out}} + d_{\text{in}} = 8{,}192$.

For a target $b \approx 0.55$: Using Eq. (17) with $b = 0.55$:

$$r = \frac{(0.55 \times 16{,}777{,}216) - 32(8{,}192)}{2(8{,}192) + 32} \approx 546$$

We choose integer $r = 546$. The actual value of $b$ for $r = 546$, using Eq. (16), is:

$$b = \frac{2(546)(8{,}192) + 32(8{,}192) + 32(546)}{16{,}777{,}216} \approx 0.5498$$

- **Example 2: Linear Layer** ($4{,}096 \times 11{,}008$) Let $d_{\text{out}} = 4{,}096$ and $d_{\text{in}} = 11{,}008$. The original number of parameters is $4{,}096 \times 11{,}008 = 45{,}090{,}816$. $d_{\text{out}} + d_{\text{in}} = 15{,}104$.
  For a target $b \approx 0.1$: Using Eq. (17) with $b = 0.1$:

$$r = \frac{(0.1 \times 45{,}090{,}816) - 32(15{,}104)}{2(15{,}104) + 32} \approx 133$$

We choose integer $r = 133$. The actual value of $b$ for $r = 133$, using Eq. (16), is:

$$b = \frac{2(133)(15{,}104) + 32(15{,}104) + 32(133)}{45{,}090{,}816} \approx 0.0999$$

By following this procedure for each linear layer in the model, we can select appropriate ranks ($r$) to achieve a desired overall value for $b$ (average bits per weight), thereby controlling the trade-off between model compression and performance. The specific ranks used to achieve the average bits per weight (BPW) figures reported in the main paper (*e.g.*, in Sections 4.1 and 5) are determined using this calculation methodology for $b$ for each layer type within the respective models.

## A.5 Pruning versus SVD

To select a base compression method to achieve sub-0.5 BPW, we compared pruning (Wanda [22]) with SVD-based compression (SVD-LLMv2 [21]) on Llama-7B ( Table 9). SVD showed superior robustness at extreme compression ratios (*e.g.*, 25% parameter retention).

Consequently, low-rank factorization was chosen to design LittleBit due to its performance resilience and the deployment advantages of dense factorized matrices over sparse pruned models (which may **require specialized hardware** or overhead). Although SVD-LLMv2 itself performs well, we found that initializing LittleBit's Dual-SVID with vanilla SVD was sufficient. The subsequent QAT process effectively recovered performance, rendering the added complexity of SVD-LLMv2 for initialization unnecessary. Accordingly, we employ the simpler vanilla SVD within our Dual-SVID initialization procedure.

Table 9: Perplexity (WikiText-2) comparison between Pruning (Wanda) and SVD-based compression (SVD-LLMv2) on Llama-7B at different parameter retention ratios. Lower perplexity is better. FP16 PPL is 5.68.

| Remaining Parameter Ratio | Pruning [22] | SVD [21] |
|---|---|---|
| 50% | 8.7 | 13.6 |
| 37.5% | 46.0 | 20.1 |
| 25% | 1842.7 | 52.6 |

## A.6 Analysis of Generated Samples

The below generated samples illustrate qualitative differences as the BPW of the LittleBit model decreases. Two main observations are:

1. Impact on Specificity, Factual Detail, and Coherence: As the effective bits decreases from 0.8 BPW towards 0.1 BPW, the generated text exhibits significant changes in detail, factual accuracy, and coherence. At 0.8 BPW, the Mona Lisa sample, while attempting to provide details (*e.g.*, "complex facial features"), becomes verbose and repetitive, and anachronistically refers to Da Vinci and the Mona Lisa in the context of the "20th century." The Turing sample at this effective bit, however, offers a standard, factually sound definition.

As the effective bit drops to 0.55 BPW, the Mona Lisa sample loses coherence significantly, expressing uncertainty about the painting's location and bizarrely shifting to a help-seeking persona ("Could you help me find the Mona Lisa painting..."). The Turing sample at 0.55 BPW begins to broaden its definition, losing some precision but remaining somewhat related.

2. Degradation in Factual Accuracy and Relevance at Extreme Compression: Extreme compression, especially at 0.1 BPW, leads to severe degradation in factual accuracy and relevance, as vividly demonstrated by both samples. The Mona Lisa sample at 0.1 BPW starts with a generally acceptable statement (*e.g.*, "iconic masterpiece") but quickly devolves into hallucination, referring to the Mona Lisa as a "statue" and associating it with "Milan's fashion world." Similarly, the Turing sample at 0.1 BPW deviates drastically from the core definition of computer science, incorrectly defining it as "the study of the theory and methods of how the human mind and the brain work" and "the science of mind." These examples underscore that while LittleBit might maintain superficial grammatical structure at very low bits, the ability to recall and generate accurate, relevant factual information is compromised, often leading to the generation of plausible-sounding but nonsensical statements.

**Sample 1: Mona Lisa**

**Prompt:** Leonardo da Vinci painted the Mona Lisa. The Mona Lisa is

| BPW | Generated Output (Phi-4 Model) |
| --- | --- |
| 0.80 | Leonardo da Vinci painted the Mona Lisa. The Mona Lisa is famous for its **complex facial features and the artist's depiction of her complex expression. It is considered one of the most important paintings of the 20th century.** The Mona Lisa is a painting by the Italian artist Leonardo da Vinci, who was a painter and sculptor. He is best known for his realistic and expressive style of painting, which is often compared to that of the natural world around him. **Leonardo da Vinci is considered one of the most important painters of the 20th century. He is famous for his complex facial features and the artist's depiction of her complex expression.** The Mona Lisa is |
| 0.55 | Leonardo da Vinci painted the Mona Lisa. The Mona Lisa is in the Louvre Museum in Paris, France. **I am not sure if it is in the museum or not. Could you help me find the Mona Lisa painting. I am searching it everywhere. It is so hard to find a picture of it. I am looking for a Mona Lisa painting with Leonardo da Vinci and the Mona Lisa by the side. Thanks in advance!** |
| 0.10 | Leonardo da Vinci painted the Mona Lisa. The Mona Lisa is an **iconic masterpiece of Italian art**, and it is one of the most photographed and photographed works of art in the world. Its execution and interpretation are breathtaking. **This statue is also one of the most popular works in Milan's fashion world**, and is one of the most photographed of the entire collection of the Mona Lisa. |

**Sample 2: Turing**

**Prompt:** Turing is called the father of computer science. Computer science is

| BPW | Generated Output (Phi-4 Model) |
| --- | --- |
| 0.80 | Turing is called the father of computer science. Computer science is a branch of mathematics that studies the design and analysis of computational systems. It is the study of the methods of designing and implementing computational systems and the application of mathematics to these systems. The most basic concepts of computer science are the concepts of computation, algorithms, data structures, and programming. |

| BPW | Generated Output (Phi-4 Model) |
|---|---|
| 0.55 | Turing is called the father of computer science. Computer science is a branch of engineering that focuses on the design, implementation, and maintenance of systems and technologies using computers. This includes designing, implementing, and maintaining the algorithms that underpin the software, hardware, and communication systems that are used in information processing, **and the Internet. Computer science is a multidisciplinary field that draws on engineering, mathematics, and business. Computer science is often divided into four broad areas** |
| 0.10 | Turing is called the father of computer science. Computer science is the study of the theory and methods of how the human mind and the brain work. **It is also called the science of mind. This is a very complex field, and many people are just not sure what to think. Some people think that computer science is the most important and most advanced field in the world, and that it is the most important branch of science in the world. This is just a myth. But there is a reason for this.** |

## A.7 Hyperparameters

Across all QAT experiments, models were trained for 5 epochs using the Adam optimizer with a cosine decay schedule. The training corpus – a mixture of C4 and WikiText-2 – contains approximately 1 billion tokens when the text is tokenized with the Llama-2-7B tokenizer ($\approx$0.2 billion tokens per epoch). The raw number of training samples (*i.e.*, the amount of original text) is identical for every model. As detailed in Table 12, for each model and BPW configuration, we swept the learning rate within the range of 4.0e-5 to 2.4e-4. We then selected the learning rate that maintained numerical stability while minimizing validation perplexity. For example, OPT-1.3B used a learning rate of 8.0e-5 at 1.0 BPW and 2.0e-4 at 0.1 BPW, while larger BPW models generally adopted slightly lower values for stability. Training was conducted using four H100 GPUs for all models except QwQ-32B, which required $4 \times 8$ A100 GPUs. In this GPU configuration notation, the first number signifies the number of nodes, and the second indicates the number of GPUs per node; thus, $4 \times 8$ represents a total of 32 GPUs.

Table 12: Knowledge Distillation Training Details

| | Training Setup | OPT | Llama | | Llama2 | | Llama3 | Phi-4 | QwQ |
|---|---|---|---|---|---|---|---|---|---|
| BPW | Target | 1.3B | 7B | 13B | 7B | 13B | 8B | 14.7B | 32B |
| 1.00 | Learning Rate | 8.0e-5 | 4.0e-5 | 4.0e-5 | 4.0e-5 | 4.0e-5 | 4.0e-5 | 4.0e-5 | 4.0e-5 |
| | # GPUs | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $4 \times 8$ |
| 0.80 | Learning Rate | 1.2e-4 | 4.0e-5 | 4.0e-5 | 4.0e-5 | 4.0e-5 | 4.0e-5 | 4.0e-5 | 1.0e-4 |
| | # GPUs | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $4 \times 8$ |
| 0.70 | Learning Rate | 1.2e-4 | 8.0e-5 | 4.0e-5 | 4.0e-4 | 4.0e-5 | 4.0e-5 | 4.0e-5 | 4.0e-5 |
| | # GPUs | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $4 \times 8$ |
| 0.55 | Learning Rate | 1.2e-4 | 8.0e-5 | 8.0e-5 | 4.0e-4 | 4.0e-5 | 8.0e-5 | 8.0e-5 | 1.0e-4 |
| | # GPUs | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $4 \times 8$ |
| 0.30 | Learning Rate | 2.0e-4 | 1.2e-4 | 8.0e-5 | 8.0e-5 | 8.0e-5 | 1.2e-4 | 8.0e-5 | 1.0e-4 |
| | # GPUs | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $4 \times 8$ |
| 0.10 | Learning Rate | 2.0e-4 | 1.2e-4 | 1.2e-4 | 1.2e-4 | 4.0e-5 | 1.2e-4 | 1.2e-4 | 1.6e-4 |
| | # GPUs | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $1 \times 4$ | $4 \times 8$ |

## A.8 Inference Efficiency Analysis

This section provides a detailed account of LittleBit's inference efficiency, complementing the kernel-level latency results presented in the main text. We analyze the theoretical computational cost, present detailed kernel latency benchmarks, and report end-to-end decoding throughput.

### A.8.1 Theoretical Cost Analysis

To complement the empirical latency results, this section provides a theoretical analysis of LittleBit's computational cost. The core principle behind LittleBit's efficiency is the replacement of a large number of expensive FP16 multiply-accumulate (MAC) operations with a smaller volume of cheaper FP16 additions and highly efficient bitwise operations (BOPs). To illustrate this, we analyze the operational cost for a single forward pass through a Llama2-7B MLP layer, where $d_{\text{in}} = 11{,}008$ and $d_{\text{out}} = 4{,}096$, quantized to 0.3 BPW, which corresponds to a latent rank of $r = 431$.

A standard FP16 matrix-vector multiplication for this layer involves $d_{\text{in}} \times d_{\text{out}}$ MAC operations. Counting each MAC as two floating-point operations (1 multiplication, 1 addition), the total computational cost for the FP16 baseline is $2 \times d_{\text{in}} \times d_{\text{out}} \approx 90.2$ million FLOPs. In contrast, LittleBit's computation is divided into floating-point and bitwise components. The FLOPs arise primarily from FP16 additions during accumulation and scaling multiplications, totaling approximately $2 \times (d_{\text{in}} + d_{\text{out}}) \times r \approx 13.0$ million FLOPs. The BOPs stem from multiplications with the binary weights ($\pm 1$), which are executed as sign-bit XOR operations, amounting to approximately $2 \times r \times (d_{\text{in}} + d_{\text{out}}) \approx 13.0$ million BOPs. This analysis reveals a nearly $7\times$ reduction in expensive FLOPs, which are replaced by an equivalent number of BOPs. As bitwise operations are substantially faster than floating-point operations on modern hardware, this theoretical breakdown provides a strong justification for the empirical latency improvements detailed below.

### A.8.2 Kernel-Level Latency

**Custom Kernel Implementation**  To empirically assess the latency of LittleBit's factorized linear layers (Eq. (5)), we developed a custom CUDA kernel. The kernel implements the two-stage computation involving the binary matrices $\mathbf{V}_{\text{sign}}$ and $\mathbf{U}_{\text{sign}}$. Key implementation details for performance include:

- **Bit-level Parallelism:** Binary weights ($\pm 1$) are packed into `uint32_t` data types. Multiplication by these weights is efficiently realized by directly manipulating the sign bit of the FP16 input values via bitwise XOR operations, avoiding costly floating-point multiplications.
- **Warp-Level Reductions:** To accelerate the accumulation step, the kernel employs warp-level reduction primitives (`warpReduceSum`). This technique efficiently sums partial results within a CUDA warp (a group of 32 threads), significantly reducing memory traffic and latency compared to naive reduction methods.
- **FP16 Arithmetic:** The scaling factors ($\mathbf{g}, \ell, \mathbf{h}$) are maintained in FP16 precision, and their application leverages native half-precision Arithmetic Logic Units (ALUs).

All latency benchmarks were performed on an NVIDIA A100 GPU with a batch size of 1. The primary baseline for comparison is a standard FP16 GEMM operation as implemented by `torch.matmul`, which leverages highly optimized libraries like CUBLAS.

**Latency Results and Discussion**  Table 13 presents the detailed latency comparison. The results demonstrate that LittleBit, accelerated by our custom kernel, can offer significant inference acceleration. For instance, in a Llama2-70B MLP-like layer, LittleBit achieves up to an **11.6× speedup** relative to the FP16 baseline at an effective BPW of 0.1. While our custom kernel is a proof-of-concept and not as exhaustively optimized as mature libraries like CUBLAS, these findings are promising. They show a clear path to substantial latency reduction at ultra-low bits, confirming that LittleBit's architecture is well-suited for high-performance deployment.

### A.8.3 End-to-End Decoding Throughput

To assess real-world application performance, we benchmarked the end-to-end decoding speed of a Llama2-7B model, measured in tokens per second (TPS). It is important to note that these throughput gains are achieved by accelerating *only* the `nn.Linear` layers with our custom kernel; other modules, such as attention and layer normalization, remained in their original FP16 implementations. Despite this partial acceleration, the results shown in Table 14 demonstrate a substantial impact on overall performance. When generating 128 new tokens, the 0.1 BPW model achieves **203.20 TPS**, a **2.46× speedup** over the FP16 baseline. This confirms that the significant kernel-level efficiencies of LittleBit translate into tangible end-to-end acceleration, even when other parts of the model leverage standard implementations.

Table 13: Kernel-level latency (ms) on an NVIDIA A100 GPU (Batch Size = 1). Compares PyTorch's FP16 GEMM with the LittleBit (LB) kernel and a baseline 1-bit (OB) kernel. Dimensions (N, M, R) denote output, input, and latent features, respectively.

| Layer Type (Model Ref.) | Dimensions (N, M, R) | Method | BPW | Latency (ms) | Relative Speedup |
|---|---|---|---|---|---|
| Llama2-70B MLP-like | $(8{,}192, 28{,}672, r)$ | FP16 Baseline | 16.0 | 0.2882 | 1.00× |
| | | OneBit | 1.0 | 0.0713 | 4.04× |
| | | LittleBit | 1.0 ($r = 6{,}400$) | 0.0938 | 3.07× |
| | | LittleBit | 0.8 ($r = 5{,}120$) | 0.0734 | 3.93× |
| | | LittleBit | 0.7 ($r = 4{,}480$) | 0.0648 | 4.45× |
| | | LittleBit | 0.55 ($r = 3{,}520$) | 0.0555 | 5.19× |
| | | LittleBit | 0.3 ($r = 1{,}920$) | 0.0372 | 7.75× |
| | | LittleBit | 0.1 ($r = 640$) | **0.0249** | **11.57×** |
| Llama2-70B ATTN-like | $(8{,}192, 8{,}192, r)$ | FP16 Baseline | 16.0 | 0.0896 | 1.00× |
| | | OneBit | 1.0 | 0.0285 | 3.14× |
| | | LittleBit | 1.0 ($r = 4{,}096$) | 0.0440 | 2.04× |
| | | LittleBit | 0.8 ($r = 3{,}296$) | 0.0340 | 2.64× |
| | | LittleBit | 0.7 ($r = 2{,}880$) | 0.0330 | 2.72× |
| | | LittleBit | 0.55 ($r = 2{,}272$) | 0.0302 | 2.97× |
| | | LittleBit | 0.3 ($r = 1{,}248$) | 0.0238 | 3.76× |
| | | LittleBit | 0.1 ($r = 416$) | **0.0207** | **4.33×** |
| Llama2-7B MLP-like | $(4{,}096, 11{,}008, r)$ | FP16 Baseline | 16.0 | 0.0620 | 1.00× |
| | | OneBit | 1.0 | 0.0226 | 2.74× |
| | | LittleBit | 1.0 ($r = 3{,}008$) | 0.0238 | 2.61× |
| | | LittleBit | 0.8 ($r = 2{,}400$) | 0.0220 | 2.82× |
| | | LittleBit | 0.7 ($r = 2{,}112$) | 0.0211 | 2.94× |
| | | LittleBit | 0.55 ($r = 1{,}664$) | 0.0192 | 3.23× |
| | | LittleBit | 0.3 ($r = 896$) | 0.0192 | 3.23× |
| | | LittleBit | 0.1 ($r = 320$) | **0.0190** | **3.26×** |

Table 14: End-to-end decoding throughput (tokens/sec) for Llama2-7B on an NVIDIA A100 GPU. The benchmark was run 10 times and averaged.

| Method | BPW | 128 New Tokens | | 256 New Tokens | |
|---|---|---|---|---|---|
| | | Avg. TPS (tok/s) | Speedup | Avg. TPS (tok/s) | Speedup |
| FP16 Baseline | 16.0 | 82.56 | 1.00× | 78.16 | 1.00× |
| LittleBit | 1.0 | 151.37 | 1.83× | 139.70 | 1.79× |
| LittleBit | 0.8 | 160.43 | 1.94× | 147.42 | 1.89× |
| LittleBit | 0.55 | 174.24 | 2.11× | 160.67 | 2.06× |
| LittleBit | 0.3 | 190.43 | 2.31× | 174.82 | 2.24× |
| LittleBit | 0.1 | **203.20** | **2.46×** | **185.39** | **2.37×** |

# B  Limitation

Despite the promising results of LittleBit in achieving extreme compression, several limitations warrant discussion. Firstly, the current LittleBit method primarily focuses on compressing the parameters within the Transformer blocks. The language model head ($lm\_head$), which typically consists of a linear layer projecting to the vocabulary size, is not subjected to the same aggressive factorization and binarization. At ultra-low bits, such as 0.1 BPW for the Transformer blocks, the $lm\_head$ can become a significant bottleneck in terms of overall model size, especially for models with large vocabularies. Thus, developing specialized compression techniques for the $lm\_head$ that are compatible with LittleBit's low-rank and binary principles is an important avenue for future work to fully realize the potential of extreme model quantization.

Secondly, while quantization-aware training (QAT) is crucial for maintaining performance at such low bits, it is computationally intensive and can be challenging to scale to extremely large models (*e.g.*, 70B parameters and beyond). Our own experiments faced resource constraints when attempting QAT for models of this magnitude. Exploring more resource-efficient QAT strategies or investigating post-training quantization (PTQ) approaches that can effectively adapt LittleBit's factorized structure to these massive models would enhance the practical applicability of our method in real-world scenarios with limited computational budgets.

Finally, the remarkable ability of models quantized with LittleBit to perform complex tasks even when retaining only a tiny fraction of the original weight information (*e.g.*, at 0.1 BPW, which can correspond to less than 1% of the original parameters' information) calls for deeper investigation.

While our empirical results demonstrate efficacy, a more fundamental understanding of how such aggressively compressed models retain their capabilities, perhaps through the lens of information theory or by analyzing changes in learned representations, would be valuable. This could lead to even more effective extreme compression techniques in the future.

## C   Societal Impact

Our work on LittleBit enhances the accessibility of large language models (LLMs) by reducing their computational and memory costs, which can foster positive societal impacts in research, education, and privacy-preserving on-device applications. However, we acknowledge that easier deployment of capable LLMs may also inadvertently lower barriers to potential misuse, such as the spread of disinformation or the amplification of societal biases. Our research responsibly utilizes existing, often publicly available models in accordance with their licenses and ethical guidelines. We strongly advocate for the ethical development and deployment of any models compressed using LittleBit, emphasizing the critical need for safeguards such as content filtering, bias mitigation techniques, transparency regarding AI-generated content, and the continuous development of robust detection mechanisms.