# Hyper-Compression: Model Compression via Hyperfunction

Feng-Lei Fan, Juntong Fan, Dayang Wang, Jingbo Zhang, Zelin Dong, Shijun Zhang, Ge Wang, Tieyong Zeng

E-mail: zeng@math.cuhk.edu.hk

September 4, 2024

## Abstract

The rapid growth of large models' size has far outpaced that of GPU memory. To bridge this gap, inspired by the succinct relationship between genotype and phenotype, we turn the model compression problem into the issue of parameter representation to propose the so-called hyper-compression. The hyper-compression uses a hyperfunction to represent the parameters of the target network, and notably, here the hyperfunction is designed per ergodic theory that relates to a problem: if a low-dimensional dynamic system can fill the high-dimensional space eventually. Empirically, the proposed hyper-compression enjoys the following merits: 1) **P**referable compression ratio; 2) **N**o post-hoc retraining; 3) **A**ffordable inference time; and 4) **S**hort compression time. It compresses LLaMA2-7B in an hour and achieves close-to-int4-quantization performance, without retraining and with a performance drop of less than 1%. Our work has the potential to invigorate the field of model compression, towards a harmony between the scaling law and the stagnation of hardware upgradation.

## 1 Introduction

Recently, the escalating demand for memory and computational resources by large models has presented formidable challenges for their deployment in resource-constrained environments [7], particularly with the stagnation of hardware Moore's Law [22]. Besides technical obstacles, the substantial cost associated with deployment also engenders profound ethical concerns. The dominance of major institutions in the realm of large models introduces risks pertaining to integrity, bias, security, and other critical aspects. Thus, how to serve large models well is not only a technological imperative but also a pursuit for social good.

One prevalent way to serve large models is to develop effective model compression approaches that crop the size of large models while maintaining acceptable performance levels. Currently, the landscape of model compression predominantly revolves around four fundamental algorithms: pruning, quantization, knowledge distillation, and low-rank approximation [29] whose basic frameworks have been established years ago. However, the compression efficacy of these methods is either capped or challenging to scale. For instance, even the most extreme case of quantization, represented by 1-bit integers, remains upset because the compression ratio is at best a constant.

Here, we introduce a novel and general-purpose approach that redefines model compression as a problem of parameter representation. This concept stems from the principles of hyper-networks, wherein smaller networks can generate weights for a significantly larger target network. We extend the concept of hypernets into what we term a 'hyperfunction' (HF), that a small parametric function can generate weights of a large network to advocate hyper-compression. Mathematically, we use a parametric function $W_n = h(n; \Theta)$ to encode the relationship between locations and weights of the target network, where $W_n$ is the $n$-th parameter of the target network, and $\Theta$ is the hyperparameters of $h$. Should the memory footprint of $\Theta$ be significantly less than that of $W$, the function $h$ can be stored on devices, allowing for the dynamic querying of weights during inference and thus facilitating the efficient deployment of large-scale models.

Table 1: The compression effectiveness of our method on UNet, MobileNetV3 and LlaMA2-7B. As a reference, the compression ratio of INT4 quantization is $4\times$.

| Model | File Size (GB) | Average (%) |
|---|---|---|
| LlaMA2-7B [21] | 12.50 | 65.88 |
| Sheared-LlaMA-1.3B [25] | 5.01 (2.50×) | 50.31 |
| TinyLlaMA [27] | 4.09 (3.06×) | 51.53 |
| LiteLlaMA[1] | 0.86 (14.53×) | 40.41 |
| LlaMA-7B + HF | 4.80 (2.60×) | 64.89 |
| Sheared-LlaMA-1.3B + HF | 0.98 (12.76×) | 49.76 |
| TinyLlaMA + HF | 0.78 (16.03×) | 50.65 |
| LiteLlaMA + HF | 0.39 (32.05×) | 39.72 |

| Model | File Size (MB) | Dice (%) |
|---|---|---|
| UNet [18] | 51.10 | 99.86 |
| Pruning | 20.30 (2.53×) | 96.34 |
| HF | 6.45 (7.93×) | 99.71 |
| Pruning + HF | 2.86 (17.87×) | 96.45 |

| Model | File Size (MB) | Top-1 Accuracy (%) |
|---|---|---|
| MobileNetV3 [13] | 5.95 | 74.41 |
| Pruning | 2.22 (2.68×) | 69.32 |
| HF | 1.18 (5.04×) | 73.93 |
| Pruning + HF | 0.43 (13.84×) | 68.47 |

One problem that ergodic theory investigates is if a low-dimensional dynamic system can eventually visit a high-dimensional space given sufficient time. We empirically find that the ergodic theory indicates a suitable hyperfunction and leads to a performant hyper-compression algorithm that offers several distinct advantages, succinctly summarized as **PNAS**: 1) **P**referable compression ratio; 2) **N**o post-hoc retraining; 3) **A**ffordable inference time; and 4) **S**hort compression time. More favorably, our method is compatible with other compression methods such as pruning to further escalate the compression efficacy. Lastly, we underscore that model compression via ergodic theory is a novel view whose intrinsic mechanism relates to chaotic theory, Kolmogorov theorem, Diophantine representation problem, etc., which remains unclear. Exploring the underlying compression mechanism will be our future work.

## 2 Result and Analysis

This section illustrates the efficacy of our novel compression methodology on three widely utilized large, middle, and small models: LlaMA2-7B [21], UNet [18], and MobileNetV3 [13]. Distinctively, our compression technique does not necessitate post-hoc retraining, even at high compression ratios and small performance drop. For example, we achieve up to a $7.93\times$ reduction in the model size of UNet within $1\%$ performance drop. This is a significant improvement over traditional methods that often require extensive retraining to restore efficacy. This capability to compress without much loss of performance potentially sets a new benchmark in model compression. Detailed descriptions of our algorithm and experiments are provided in **Supplementary Information**.

Table 2: Changes in inference time on UNet

| Batch Size | 1 | 2 | 4 | 8 | 16 |
|---|---|---|---|---|---|
| **Normal-HF (s)** | 5.70 | 5.88 | 6.00 | 7.36 | 15.64 |
| **Optimized-HF (s)** | 1.17 | 1.22 | 1.39 | 2.70 | 11.12 |
| **Original (s)** | 0.18 | 0.28 | 0.48 | 1.78 | 10.07 |

**Preferable compression ratio**. As shown in Table 4, we apply the same pruning method used in Sheared-LlaMA-1.3B, TinyLlaMA, and LiteLlaMA, combined with our compression technique. We evaluate our compressed models on eight downstream tasks: 0-shot accuracy on SciQ, WinoGrande, ARC-E, 25-shot ARC-C, 10-shot

---

[1]For more details, please visit the Hugging Face website of LiteLlaMA: `https://huggingface.co/ahxt/LiteLlama-460M-1T`.

Table 3: The compression time of our method on UNet, MobileNetV3
and LlaMA2-7B

| Model | Compression Time (s) |
|---|---|
| UNet + HF [18] | 30.00 |
| Pruned-UNet + HF | 34.60 |
| MobileNetV3 + HF | 11.63 |
| Pruned-MobileNetV3 + HF | 8.83 |
| **Model** | **Compression Time (min)** |
| LlaMA2-7B + HF | 53 |
| Sheared-LlaMA-1.3B + HF | 7 |
| TinyLlaMA + HF | 11 |
| LiteLlaMA + HF | 6 |

HellaSwag, 32-shot BoolQ, NQ, and 5-shot MMLU. The average test results of eight downstream tasks are shown in Table 4, where "File Size" indicates the size of the model parameter file. The comparison of individual tasks is put in **Supplementary Information**. It is noteworthy that our method can compress LlaMA2-7B by a factor of $2.60\times$ while maintaining the Average score decrease within 1%, whereas Sheared-LlaMA-1.3B achieves $2.60\times$ with $15.57\%$ performance decreases.

As for UNet and MobileNetV3, **our method can compress UNet and Pruned-UNet by** $7.93\times$ **and** $7.10\times$ **with the performance loss contained in 1%**. Particularly, our method succeeds in combination with other model compression methods such as pruning to achieve an even higher compression ratio. In UNet, **the total compression ratio is** $17.87\times$ **with the performance loss** $3.41\%$.

**No post-hoc retraining**. Despite the presence of minor discrepancies between the original and decoded parameters—specifically, at magnitudes ranging from $10^{-4}$ to $10^{-3}$—the impact on error accumulation through layer-by-layer information propagation remains acceptable. This is due to the inherent redundancy within the network itself, including parametric redundancy and activation space redundancy. As a result, the destructive effect of our method is minimal, obviating the need for post-hoc retraining. This characteristic is particularly advantageous in real-world settings, where often 1) training data are inaccessible, and curating data for fine-tuning is costly; 2) no computing resources are supplied for retraining.

**Affordable inference time**. A network is a hierarchical structure. To expedite the inference time, our approach is to parallelize the decoding and inference processes. The key is to complete the decoding of a layer before using this layer to infer. As shown in Table 2, "Normal-HF" means the model without decoding and inference parallized. "Original" refers to the time required for a single inference of the original UNet, and "Optimized" refers to the time required for one inference by employing parallel decoding and inference processes. when batch size increases, the inference time of our compressed model gets closer to that of the original model, which aligns with our expectations. This is because when batch size goes up, the inference time across layers is slow, which allows weights in later layers to be decompressed in time.

**Short compression time**. As Table 6 shows, our method can compress models very fast. This efficiency primarily stems from matrix operations to expedite the majority of the compression process and the K-D tree to expedite the search. Additionally, by treating the compression tasks between layers as independent operations, we implement a multiprocessing parallel strategy to further decrease compression times.

# 3  Materials and Methods

## Motivation

A human genome of only approximately 30,000 genes can remarkably encode the development of a human brain that has finally billions of neurons and trillions of synaptic connections [20]. This observation suggests the existence of a concise and implicit mapping from genotype to phenotype, capable of expanding a limited number of genetic instructions into a vast array of biological functionalities. Inspired by this efficient genetic encoding, the idea of hypernet is to use a small network (genotype, called hypernet) to generate weights for the target network (phenotype).

We consider the hypernet in the setting of model compression. Furthermore, we generalize the idea of hypernet to a parameterized function (hyperfunction). The biological observation is the existence of a concise
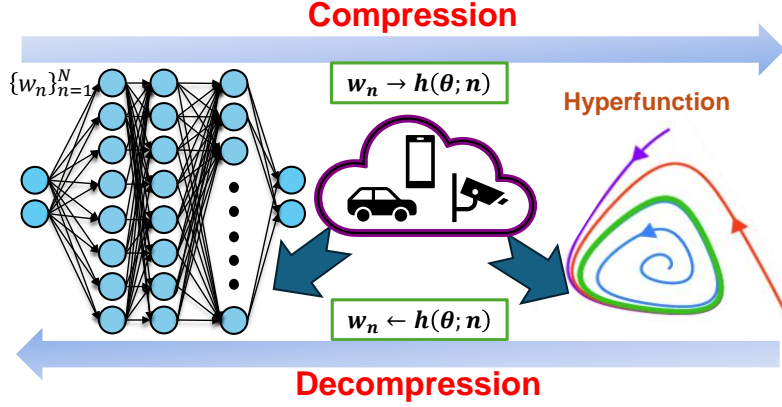
Figure 1: We use a parametric function $W_n = h(n; \Theta)$ to encode the relationship between locations and weights of the target network, where $W_n$ is the $n$-th parameter of the target network. $h$ is implied by the ergodic theory.

mapping, regardless of whether it is denoted by a network or other forms of functions. With a hyperfunction, instead of directly diminishing weights, one stores a hyperfunction than the model and queries its weights when needed. This novel perspective converts the model compression problem into the parameter representation issue. Mathematically, we use a hyperfunction $W_n = h(n; \Theta)$ to encode the relationship between locations and weights of the target network, where $W_n$ is the $n$-th parameter of the target network, and $\Theta$ is the hyperparameters of $h$. Along this direction, the important question is *how can we design a suitable hyperfunction?*

## Tools from Ergodic Theory

The above question can be affirmatively answered through the lens of ergodic theory. One branch of the ergodic theory [6] investigates if the trajectory of a low-dimensional dynamical system will eventually traverse every point of a given high-dimensional space. From the engineering point of view, the ergodic theory discusses the feasibility of using low-dimensional curves to approximate points in high-dimensional space, which indicates a potentially ideal hyperfunction for model compression.

**Theorem 1.** *For any $K \in \mathbb{N}^+$, suppose $\alpha_1, \cdots, \alpha_K$ are irrationally independent, the following point set*

$$\left\{ \left[ \tau(a_1 \cdot \theta), \tau(a_2 \cdot \theta), \cdots, \tau(a_K \cdot \theta) \right]^T : \theta \in [0, +\infty] \right\} \tag{1}$$

*is dense in $[0,1]^K$, where $\tau(x) = x - [x]$.*

In Theorem 1, the irrationality condition is easily fulfilled. For example, we can set $a_k = 1/(\pi + k)$. Given $\{W_n\}_{n=1}^N$, according to Theorem 1, we know that for any $\epsilon > 0$, there exists $\theta^*$ such that

$$|W_n - \tau(\theta^*/(\pi + n))| < \epsilon, n = 1, \cdots, N. \tag{2}$$

Theorem 1 provides a natural and elegant way to encode a group of numbers into one number. However, we should only be cautiously optimistic about it since simply stitching digits of numbers can also ensemble a group of numbers into one number. What matters is whether we can find a simple $\theta$ that approximates $W_1, \cdots, W_K$ based on (2), such that the memory footprint of the former is smaller than the latter. Our experiments empirically confirm that the parametric representation according to the ergodic theory has compression efficacy. This is primarily because our compression is lossy, and the network itself has redundancy.

## Core Issues in Algorithmic Development

The entire algorithmic development is anchored in (2). It begins by partitioning the parameters $\{W_n\}_{n=1}^N$ into groups, each with $K$ parameters. The core objective of the algorithm is, given $\alpha_1, \cdots, \alpha_K$ and an acceptable error $\epsilon$, to identify the smallest $\theta^*$ for points in $[0,1]^K$, aiming for the highest feasible dimension $K$. First, though the existence of $\theta^*$, because a network usually has millions or billions of parameters, searching $\theta^*$ for a group of

$K$ parameters has to be fast. Second, an imprecise $\theta^*$ could lead to significant performance degradation in the restored networks. Therefore, the algorithm must ensure that $\theta^*$ approximates the parameters precisely. Third, when time and precision allow, the algorithm should return the smallest $\theta^*$ for the highest compression ratio. The algorithm details are shown in **Supplementary Information**.

## 4 Conclusion and Future Work

In this study, we have proposed hyper-compression, a novel and general-purpose methodology for model compression that leverages hyperfunction to encode the parameters of the target network. We have conducted experiments on various architectures, including LlaMA2-7B, UNet, and MobileNetV3, demonstrating that our approach is both user-friendly and scalable. Nevertheless, it represents merely an initial foray into the whole blueprint. Future directions include 1) theoretically describing the relationship among the error, the length of $\theta$, and the dimension $K$ to reveal the essence of using the ergodic theory to realize compression; 2) beyond classical theory, exploring other possibilities such as implicit neural representation [3] to strike a better balance between compression time, inference time, and compression ratio. We believe that hyper-compression can contribute to Moore's law of model compression in the near future, *i.e.*, the compression efficiency of model compression algorithms can be doubled annually, to meet the need of large models's development.

# 5    Supplementary Information

# Contents

# I. Limit of the Compression Ratio of Pruning Convolutional Networks

In the main body, we mention that pruning-based methods are hard to scale. Here, we show that the compression ratio of pruning a convolutional network is constantly bounded by studying a randomly generated infinite network and introducing the percolation theory [10] that characterizes the connectivity of a network when links are removed per certain rules.

We study one-dimensional convolutional networks whose topology is shown in Figure 2. We assume any edge of a network is randomly removed with the same probability. This is also equivalent to the following process: Randomly generate weights for edges of a given network. If the weight of an edge is higher than a threshold, the edge is retained; otherwise, it is removed. In the jargon of percolation theory, an edge is open or closed if the edge exists or disappears. When the probability of removing edges reaches a threshold, referred to as the percolation threshold, there will be no connected clusters from the input and output, meaning that the network will always yield 0 for any input. The probability threshold indicates the maximal proportion of a network's edges to prune. This section discusses the upper and lower bounds of the percolation threshold for a class of graphs defined by the one-dimensional convolutional network, thereby quantitatively characterizing the compression ratio.

Let us first provide relevant definitions below:

**Definition 1** (Graph, Vertex, Edge). *A graph $G = (V, E)$ is an ordered pair of sets, where $V$ is the vertex set, and $E$ is the edge set. In the edge set $E$, we present the edge between $v, w \in V$ as $\langle v, w \rangle$.*

The class of graphs we investigate is $G_r = (V_r, E_r)$, $r \in \mathbb{Z}_{\geq 2}$, where $V_r = \mathbb{Z}^2$, $E_r = \{\langle (m, n), (m+1, n+i) \rangle \mid m, n \in \mathbb{Z}, i \in 0, 1, \cdots, r-1\}$, that is, a convolution with the kernel size $r$ and stride 1.

**Definition 2** (Open and Closed). *Let $p \in [0, 1]$. Declare each edge of $G_r$ to be open with probability $p$ and closed otherwise, independently of all other edges.*

**Definition 3** (Open Cluster). *Consider the random subgraph of $G_r$ containing the vertex set and the open edges only. The connected components of this graph are called open clusters. We write $C_r(v)$ for the open cluster containing the vertex $v$, and we call $C_r(v)$ the open cluster at $v$.*

**Definition 4** (Percolation probability and critical probability). *Let $\theta_r(p)$ being the probability that given vertex, $v$, belongs to an infinite open cluster in $G_r$, that is $\theta_r(p) = P(\mid C_r(v) \mid = \infty)$. By the translation invariance of those graphs, we lose no generality by taking this vertex as the origin and denoting $C_r(0, 0)$ as $C_r$. The critical probability is defined to be $p_{c,r} = \sup\{p \in [0, 1] \mid \theta_r(p) = 0\}$.*

Earlier, researchers have mainly focused on planar graphs, but we observe that $G_r$ is a planar graph only when $r = 2$, and is non-planar for $r \geq 3$, part of $G_2$ and $G_3$ are shown in Figures 2 and 3. Therefore, many methods, such as using the duality of planar graphs and the star-triangular transition, cannot be applied. In the following, we will first solve the critical probability for the case of $r = 2$, and then study the upper and lower bounds of the critical probability for $r \geq 3$. Our method is degenerating a non-planar graph into a planar graph by removing edges.

Figure 2: $G_2$      Figure 3: $G_3$      Figure 4: $G_{-1}$      Figure 5: $G_{triangular}$

**Proposition 1** (Critical Probability of $G_2$). *$p_{c,2} = 0.5$.*

*Proof of Proposition 1.* Let $\mathbb{L}^2$ be the square lattice on $\mathbb{Z}^2$, whose percolation threshold is well-known to be 0.5 [10]. Consider $\phi : \mathbb{Z}^2 \longrightarrow \mathbb{Z}^2$ defined by $\phi(m, n) = (m - n, n)$. We are going to show $\phi$ is graph isomorphism.

For points $(m, n), (m + 1, n)$ connected by the edge $\langle (m, n), (m + 1, n + i) \rangle$ in $G_2$, we have $\phi(m, n) = (m - n, n), \phi(m + 1, n + i) = (m - n + 1, n)$ are connected by $\langle (m - n, n), (m - n + 1, n) \rangle$ in $\mathbb{L}^2$.

For points $(m, n), (m + 1, n + 1)$ connected by the edge $\langle (m, n), (m + 1, n + 1) \rangle$ in $G_2$, we have $\phi(m, n) = (m - n, n), \phi(m + 1, n + 1) = (m - n, n + 1)$ are connected by $\langle (m - n, n), (m - n, n + 1) \rangle$ in $\mathbb{L}^2$.

Therefore, $\phi$ is a graph isomorphism from $G_2$ to $\mathbb{L}^2$, and hence the percolation threshold of $\mathbb{L}^2$ is the same as that of $G_2$, which is 0.5.

$\square$

**Proposition 2** (Critical Probability of $G_r$, $r \geq 3$). *For $r \geq 3$, $\frac{1}{\mu_r} \leq p_{c,r} \leq p_0$, where $\mu_r$ is the connected constant of $G_r$ and $p_0$ is the solution of $2p + p^2 - p^4 = 1$ between 0 and 1.*

*Proof of Proposition 2.* The proof of the part $\frac{1}{\mu_r} \leq p_{c,r}$ is a more general result that the percolation threshold of any graph is no less than the reciprocal of its connective constant, see [10, 2].

The proof of the part $p_{c,r} \leq p_0$ is even more complicated. Define a graph $G_{-1} = (V_{-1}, E_{-1})$, where $V_{-1} = \mathbb{Z}^2$, $E_{-1} = \{\langle (m, n), (m + 1, n + i) \rangle \mid m, n \in \mathbb{Z}, i \in \{-1, 0, 1, \cdots, r - 2\}\}$, and Figure 4 shows the $G_{-1}$ corresponding to $G_3$. Meanwhile, consider the function $\Phi : \mathbb{Z}^2 \longrightarrow \mathbb{Z}^2$ defined by $\phi(m, n) = (m, n - m)$, which maps the two points $(m, n), (m + 1, n + i)$ connected by an edge in $G_r$ to $(m, n - m), (m - n, n - m + i - 1)$ that are connected by an edge in $G_{-1}$. Thus, $\Phi$ is a graph isomorphism from $G_r$ to $G_{-1}$, and they have the same percolation threshold, i.e., $p_{c,r} = p_{c,-1}$.

Next, we modify $G_{-1}$ by removing or connecting vertices and edges.

Remove vertices: $(2m + 1, 2n), (2m, 2n + 1), m, n \in \mathbb{Z}$,

Remove edges: $\langle (2m+1, 2n), (2m+2, 2n+i) \rangle, \langle (2m, 2n+1), (2m+1, 2n+j), \rangle, i = -1, 1, \cdots, r-2, j = 0, 2, \cdots, r - 1$.

Connect pair of edges $\langle (2m, 2n), (2m + 1, 2n) \rangle$ and $\langle (2m + 1, 2n), (2m + 2, 2n) \rangle$ into $\langle (2m, 2n), (2m + 2, 2n) \rangle$.

Connect pair of edges $\langle (2m, 2n + 1), (2m + 1, 2n + 1) \rangle$ and $\langle (2m - 1, 2n + 1), (2m, 2n + 1) \rangle$ into $\langle (2m - 1, 2n + 1), (2m + 1, 2n + 1) \rangle$.

The resulting graph is a triangular lattice (rotated 45 degrees counterclockwise to the usual form), as shown in Figure 5, denoted as $G_{triangle}$. During the construction of $G_{triangle}$, each connected subgraph in $G_{triangle}$ can uniquely correspond to a connected subgraph in $G_{-1}$, i.e., those untouched edges and vertices are naturally embedded, those connected edges are restored to one vertex and two edges.

To maintain the correspondence with $G_{-1}$ probabilistically, the edges formed by merging in $G_{triangle}$ have probability $p^2$, while the remaining edges have probability $p$ of being open. Under the probability, the probability that each connected subgraph in $G_{triangle}$ is open is the same as the probability that its corresponding subgraph in $G_{-1}$ is open. Furthermore, each infinite open cluster in $G_{triangle}$ can uniquely correspond to an infinite open cluster in $G_{-1}$, but not vice versa. Therefore, $\theta_{-1}(p) = P(|C_{-1}(v)| = \infty) \geq P(|C_{triangle}(v)| = \infty) = \theta_{triangle}(p)$, and thus $p_{c,-1} \leq p_0$.

Fortunately, this inhomogeneous triangular lattice $G_{triangle}$ has been studied extensively [10], and the solution under the provided three-edge probabilities is the solution of $2p + p^2 - p^4 = 1$ between 0 and 1, which is approximate 0.425787.

$\square$

**Remark**. From the viewpoint of percolation theory, the compression ratio of pruning a randomly generated convolutional network is $1/p_0 < \frac{1}{p_{c,r}} < \mu_r$, which is a fixed constant given $r$. Our analysis is based on an assumption that ignores the association between two edges. In reality, when an edge with a larger weight is pruned, edges with smaller weights are also pruned; otherwise, they are not. Future work will consider the association to have a more accurate estimation.

# II. The Trend of Growth of Large Models

In recent years, the field of natural language processing (NLP) has witnessed remarkable advancements, particularly with the development of large language models (LLMs). These large models have demonstrated impressive capabilities in various language-related tasks, with a rapid increase in size over the past few years. For instance, GPT-3 consists of a staggering 175 billion parameters. Recently, M6, one of the largest models to date, has over $10^4$ billion parameters. This growth in model size has been driven by the desire to improve the model's complex language understanding by scaling law. However, the growth of GPU memory has not kept pace. GPUs, which are crucial for training and serving these models, have seen relatively modest increases in memory capacity. This disparity between models' size and GPU memory has created significant challenges in the development of large models.
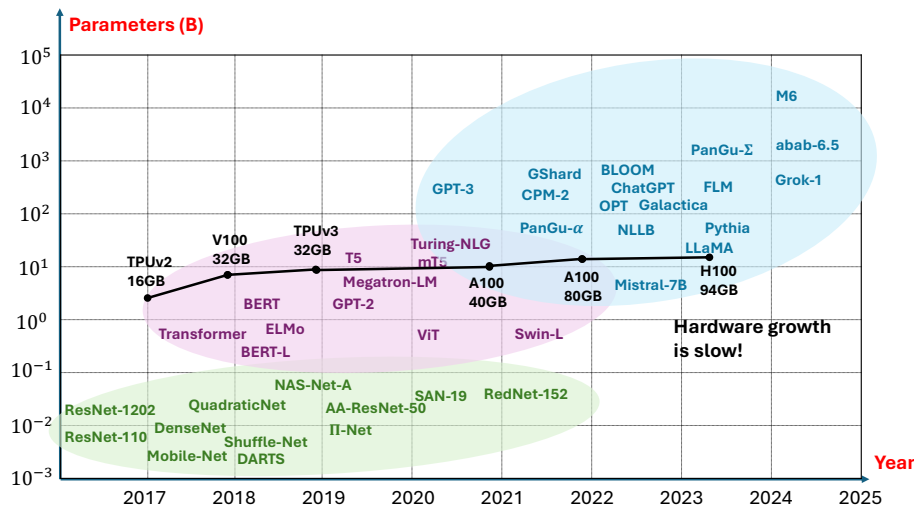


Figure 6: The rapid growth of large models' sizes has outpaced the growth of GPU memory, creating huge challenges in serving large models.

Specifically, the disparity between the model size and GPU memory brings the following problems:

1. **Memory Bottleneck**: The limited GPU memory poses a bottleneck for serving LLMs. When deploying LLMs for real-world applications, the entire model and its associated parameters need to fit within the GPU memory. However, as LLMs continue to grow in size, it becomes increasingly difficult to accommodate them within the available memory constraints. This limitation hampers the deployment and wide adoption of LLMs in practical scenarios.

2. **Ethic Issues**: Since large models are hard to serve, large institutes will gradually dominate the development and deployment of large language models, which raises ethical concerns:

   Algorithmic Bias: LLMs developed by dominant institutes may reflect the biases present in the data they were trained on. This can perpetuate existing societal biases and discrimination, leading to unfair outcomes in decision-making processes.

   Lack of Transparency: Large institutes may not always disclose the full details of their model architectures, training data, or decision-making processes. This lack of transparency can hinder accountability, making it difficult to assess the fairness and reliability of LLMs in real-world applications.

   Ethical Use Cases: The ethical implications of deploying LLMs in various applications, such as content moderation, surveillance, or healthcare, need to be carefully considered. Large institutes' control over these technologies raises concerns about how they are used and the potential impact on individuals and society.

As the demand for more powerful models and the risk of monopoly by large institutes continue to rise, addressing the disparity between model size and GPU memory is a paramount research topic. Particularly, we need

to cultivate novel model compression techniques that can realize Moore's law of model compression, facilitating the deployment and efficient utilization of LLMs in practical scenarios. Moore's law of model compression, which proposes that the compression ratios of models should double annually, is the potential solution. By compressing models, we can reduce their size and memory requirements, making them more feasible to deploy on GPUs with limited memory.

# III. Algorithmic Details

Anchored in the hyperfunction driven by ergodic theory mentioned in the main body, we design a general-purpose model compression method that enjoys the following merits (**PNAS**): **P**referable compression ratio, **N**o post-hoc retraining, **A**ffordable inference time, and **S**hort compression time. In this section, we describe in detail our entire model compression method which is made of four parts, and each part corresponds to one facet of merits.
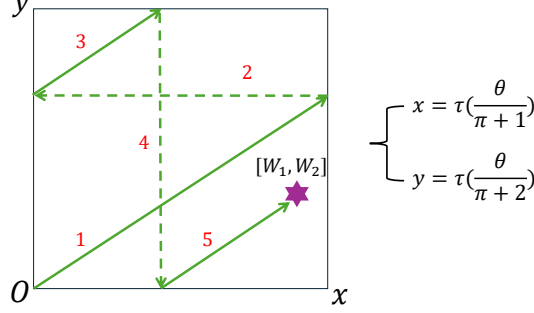


Figure 7: An illustrative scheme for using ergodic theory to learn a low-dimensional representation of high-dimensional points.

**Preferable compression ratio**

The core objective of the algorithm is to identify the smallest $\theta^*$ for points in hypercubes $[0,1]^K$ under the condition that $\epsilon$ is sufficiently small and aiming for the highest feasible dimension $K$. To this end, we consider the distribution of weights of the original network by binning weights into three groups based on magnitudes of weight values.

- First, we flatten all parameters into a one-dimensional vector $\mathbf{W} = [W_1, W_2, \cdots, W_N]$, where $N$ is the total number of weights, and then split it into groups $[\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \cdots, \mathbf{W}^{(G)}]$, where $W^{(g)} = [W_{K(g-1)+1}, W_{K(g-1)+2}, \cdots, W_{K(g-1)+K}], g = 1, 2, \cdots, G$ is of $\mathbb{R}^{1 \times K}$, and $G$ is the number of groups. For any $\mathbf{W}^{(g)}$, there exists an integer $\theta_g^*$ such that $|\mathbf{W}^{(g)} - \tau(\theta_g^* \cdot \mathbf{a})| < \epsilon, g = 1, \cdots, G$ and $\mathbf{a} = [a_1, a_2, \cdots, a_K]$. Finally, the vector $[\theta_1^*, \theta_2^*, \cdots, \theta_G^*]$ represents a compression to $\mathbf{W}$. Take $K = 2$ and $\mathbf{W} = [W_1, W_2, W_3, W_4] = [0.1, 0.2, 0.4, 0.5]$ as an example, we split $\mathbf{W}$ as $[\mathbf{W}^{(1)}, \mathbf{W}^{(2)}]$, where $\mathbf{W}^{(1)} = [W_1, W_2]$ and $\mathbf{W}^{(2)} = [W_3, W_4]$, and derive $\theta_1^*$ and $\theta_2^*$ for $\mathbf{W}^{(1)}$ and $\mathbf{W}^{(2)}$, respectively. As Figure 7 shows, we can find a $\theta$ for $(x, y) = (\tau(\theta/(\pi + 1)), \tau(\theta/(\pi + 2)))$ that approximates a point in 2D space.

- Given a parameter vector $\mathbf{W}$, $[\theta_1^*, \theta_2^*, \cdots, \theta_G^*]$ is a compression. However, we do not want $\max(\theta_i^*)$ to be a large number, as this would require more storage space. Therefore, we define an integer $U$ as the upper bound of $\theta_g^*$ to ensure the preferable compression efficacy. Next, $[\theta_1^*, \theta_2^*, \cdots, \theta_G^*]$ can be stored using the data type uint$m$, where $m = \lfloor \log_2(\max(\theta_i^*)) \rfloor + 1$, allowing $[\theta_1^*, \theta_2^*, \cdots, \theta_G^*]$ to be stored using the minimum possible number of bits without causing data overflow. In other words, given a target two-dimensional point $\mathbf{W}^{(g)}$ and $U$, what we are doing is finding an integer $\theta_i^* \in [0, U]$ such that the Euclidean distance between $\tau(\theta_g^* \cdot \mathbf{a})$ and $\mathbf{W}^{(g)}$ is minimized. As $U$ increases, potentially, more integers can be explored to make the loss $|\mathbf{W}^{(g)} - \tau(\theta_g^* \cdot \mathbf{a})|$ smaller. Each layer can select a different $U$, allowing important layers to choose a larger $U$ to ensure the approximation error remains sufficiently small.

- In our algorithm, we set $K = 2$ and replace the unit square $[0, 1]^2$ with a more flexible box $[a, b] \times [c, d]$. Specifically, given $\mathbf{W}$, $[a, b] \times [c, d]$ is defined as a square with side length $l$, centered at $(\bar{x}_i, \bar{y}_i)$, where $(\bar{x}_i, \bar{y}_i)$ is the centroid of two-dimensional points $\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \cdots, \mathbf{W}^{(G)}$. $l$ is a hyperparameter and typically set to 0.01 for optimal results. Consequently, the function $\tau(x) := x - \lfloor x \rfloor$ in the ergodic theorem is redefined as $\tau(z) := f(z) + ([\bar{x}_i, \bar{y}_i] - [\frac{l}{2}, \frac{l}{2}])$, where $f(z) = z \bmod l$.

- In the ergodic theorem, $\mathbf{a} = [a_1, a_2] = \mathbf{d} \cdot I$ is a two-dimensional directional vector, where $\mathbf{d}$ and $I$ determine the direction vector and step size for the movement of curves defined by the parametric equation, respectively. Since $\theta$ is an integer no more than $U$, at most we have $U + 1$ points in two-dimensional space $(\tau(a_1\theta), \tau(a_2\theta)), \theta = 0, 1, \cdots, U$. In our experiment, we define an optimal vector $\mathbf{a}$ based on $U$, such that the distribution of $U + 1$ points is more uniform, thereby minimizing the maximum error as much as possible. Specifically,

$$
\begin{aligned}
\mathbf{a} &= \mathbf{d} \cdot I \\
&= [\frac{l}{U}, l] / \|\frac{l}{U}, l\| \cdot I \\
&= [\frac{l}{U}, l] / \|\frac{l}{U}, l\| \cdot \frac{l}{\sin \alpha \cdot \lfloor \sqrt{U} \rfloor},
\end{aligned}
\tag{3}
$$

where $\tan \alpha = \frac{l}{l/U} = U$. As shown in Fig. 8, it demonstrates that using (3) results in a more uniform distribution of the sample nodes compared to defining $\mathbf{a}$ as follows:

$$
\mathbf{a} = [1/(\pi + 1), 1/(\pi + 2)].
\tag{4}
$$



Figure 8: Different $\mathbf{a}$ will result in different distribution of points. Upper three subfigures use (3), while lower three figures utilize (4).

- If we construct a square with its center at the centroid $(\bar{x}_i, \bar{y}_i)$ and side length $l$, it is highly probable that many points in $\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \cdots, \mathbf{W}^{(G)}$ will fall outside this square. Therefore, we need to first scale these points into the square using a scaling factor $s$ to compress, and then scale the substituted points $\tau(\theta_g^* \cdot \mathbf{a})$ back during decompression. While the error rate keeps intact, this inversion will amplify the error $\epsilon$, where $\epsilon = |W^{(g)} - \tau(\theta_g^* \cdot a)|$. It should be noted that, assuming $\mathbf{W}^{(F)}$ is the point farthest from the centroid $(\bar{x}_i, \bar{y}_i)$ in $\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \cdots, \mathbf{W}^{(G)}$, we define the farthest distance $l_f$ as $|\mathbf{W}^{(F)} - [\bar{x}_i, \bar{y}_i]|$, and then the scaling factor $s = \frac{l}{2 \cdot l_f}$. The specific formula derivation is as follows:

$$
\begin{aligned}
\mathbf{W}^{(in)} &= \mathbf{W}^{(in)} - C \\
&= C\mathbf{W}^{(out)} \cdot s \\
&= C\mathbf{W}^{(out)} \cdot \frac{l}{2 \cdot l_f},
\end{aligned}
\tag{5}
$$

where $C$ is the centroid node, $\mathbf{W}^{(out)}$ is a node ouside the squre, $\mathbf{W}^{(in)}$ is the scaled node of $\mathbf{W}^{(out)}$.

We denote $\tau(\theta_i^* \cdot \mathbf{a})$ as $\mathbf{W}^{(in)*}$, so that $\epsilon = \mathbf{W}^{(in)*} - \mathbf{W}^{(in)}$. The detailed scale-back process is shown as follows:

$$\begin{aligned}
\mathbf{W}^{(out)*} &= \frac{C\mathbf{W}^{(in)*}}{s} + C \\
&= \frac{\mathbf{W}^{(in)} + \epsilon}{s} + C
\end{aligned} \tag{6}$$

Therefore, we can derive the error:

$$\begin{aligned}
error = |\mathbf{W}^{(out)} - \mathbf{W}^{(out)*}| &= |C\mathbf{W}^{(out)} - C\mathbf{W}^{(out)*}| \\
&= |C\mathbf{W}^{(out)} - \frac{C\mathbf{W}^{(in)*}}{s}| = |\frac{\mathbf{W}^{(in)*} - \mathbf{W}^{(in)}}{s}| \\
&= |\frac{\epsilon}{s}| = |\frac{2 \cdot \epsilon \cdot l_f}{l}|
\end{aligned} \tag{7}$$

From (7), we can observe that the farthest distance $l_f$ determines the error of estimating $\mathbf{W}^{(out)}$ after scaling back.



Figure 9: Given K=2, the points outside the center square are divided into three categories: blue points belong to the first category, pink points belong to the second category, and green points belong to the third category.

Therefore, as shown in Fig. 9 in our algorithm, we classify the points outside the square into $M$ categories based on their distances from the centroid $C$, where $M$ is a hyperparameter, so that the points in different categories can utilize different farthest distance $l_{fm}$ to define different scale factor $s_m$. This can mitigate the adverse impact of the global farthest distance $l_f$ on the final error as much as possible. This method is highly effective in practice, as a significant portion of the parameter nodes $\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \cdots, \mathbf{W}^{(G)}$ exhibit a pattern that most points are close to the centroid.

Specifically, assuming the points outside the square are divided into $M$ categories and $W^{(m)}$ is a point that belongs to the $m$-th category, the scaling factor $s_m$ for this point is defined as

$$s_m = \frac{l/2}{l/2 + (l_f/K) \cdot m}. \tag{8}$$

The final compression $\theta_m^*$ is given by the following formula:

$$\theta_m^* = m \cdot U + \lambda_m^*, \tag{9}$$

where $\lambda_m^*$ is the compression of the scaled $W^{(m)}$ in the center square.

**No fine-tuning**

Our method does not require fine-tuning of the model, and can directly use the trained model for compression. This greatly simplifies the model compression process and avoids the time and computing resource consumption caused by retraining the model. The reasons why our algorithm requires no retraining are mainly because our algorithm can guarantee the error rate is low by optimizing the direction of $\mathbf{a}$, the side length $l$, and $U$. The error rate of points in different classes is around $(\frac{l}{U})/l = \frac{1}{U}$.

**Affordable inference time**

At first glance, one may think that model compression based on hyperfunctions suffers from a slow inference time, as this technique needs to restore parameters before inference, which adds another level of computation. Here, we propose a scheme that leverages the intrinsic hierarchical structure of a network to greatly reduce the inference time. Our scheme parallelizes parameter decompression and inference. As shown in Figure 10, while the parameters of later layers (except the first part) are restored, the inference operation in earlier layers is also carried out simultaneously. As long as we can recover the parameters of the current layer before the inference arrives at the current layer, there is no waste of time. Thus, theoretically, the inference time of using our algorithm only increases quite moderately. The increment is the time used to restore the parameters of the first layer.



Figure 10: The rapid growth of LLM sizes has outpaced the growth of GPU memory, creating challenges in serving these increasingly massive models.

Based on the above solution, we need to shorten the decompression time. First, the entire decompression process is implemented using matrix operations, which can significantly speed up the computation. Second, the process of reading files from storage and preprocessing them is relatively time-consuming, and repeats each time the function is called. Therefore, we optimize this process, performing file reading and preprocessing operations only when the function is called for the first time, and storing the processing results in the cache. When the function is called again, the preprocessed intermediate results are directly captured from the cache, further shortening the inference time.

**Short Compression Time**

With the advent of large models, compression time becomes an important facet of evaluating a compression algorithm. Here, we propose a suite of techniques to enable the proposed hyper-compression method to fast complete the model compression process with off-the-shelf tools, for both CPU and GPU inference.

- As mentioned earlier, when encoding the parameter vector $\mathbf{W}$, we first convert it into many points $\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \cdots, \mathbf{W}^{(G)}$. Then, given $\mathbf{W}^{(g)}$, we search $u$ from $\tau(1 \cdot \mathbf{a}), \cdots, \tau(U \cdot \mathbf{a})$ to approximate

$\mathbf{W}^{(g)}$. Next, instead of repeating searching for every point, we store $\tau(1 \cdot \mathbf{a}), \cdots, \tau(U \cdot \mathbf{a})$ in the format of k-d trees [17]. Thus, we can fast determine $\tau(u \cdot \mathbf{a}) \in [\tau(1 \cdot \mathbf{a}), \cdots, \tau(U \cdot \mathbf{a})]$ that has the shortest distance to $\mathbf{W}^{(g)}$.

- We extensively use matrix operations. Given a series of two-dimensional point coordinates $\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \cdots, \mathbf{W}^{(G)}$, we first calculate the scaling factor list $F$:

$$F = [s_1, s_2, \cdots, s_G]^\top,$$
$$s_g = \frac{l/2}{l/2 + (l_f/M) \cdot m_g}, \tag{10}$$

where $s_g \in F$ represents the scale factor of $\mathbf{W}^{(g)}$, $m_g$ means the category of node $\mathbf{W}^{(g)}$ and the points located within the center square are classified as category $m_g$.

Let $\mathbf{m} = [m_1, y_2, \cdots, m_G]^\top$, $\mathbf{W} = \left[\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \cdots, \mathbf{W}^{(G)}\right]^\top$. Then, we can calcute the final compression $\theta^*$ as follows:

$$C\mathbf{W} = \mathbf{W}^\top - [C, \cdots, C]^\top,$$
$$O = C\mathbf{W} \cdot F + [C, \cdots, C]^\top$$
$$\lambda^* = KD(O^\top, [\tau(1 \cdot a), \cdots, \tau(U \cdot \mathbf{a})]) \tag{11}$$
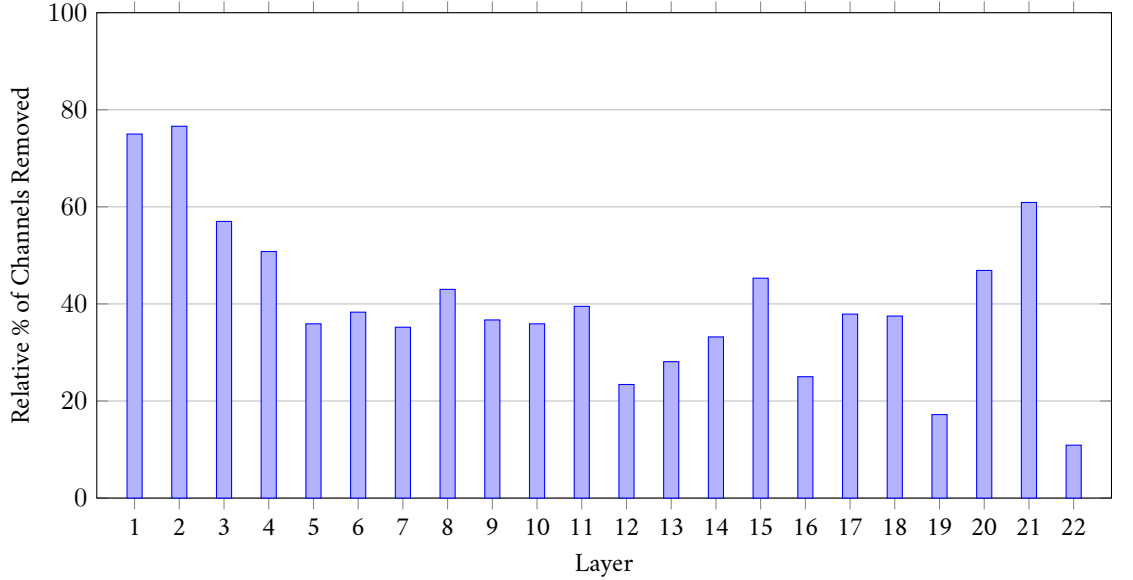$$\theta^{*\top} = \lambda^{*\top} + U \cdot \mathbf{m}$$

- We also multiprocess the compression by simultaneously compressing different layers of a model to maximize the utilization of computational resources.

# VI. Experimental Details

In this section, we introduce details of experimental configurations. All experiments are conducted on a single NVIDIA RTX4090 (24GB) GPU using the PyTorch framework.

- **Details of experiments related to UNet**: We perform parameter compression on both UNet and the pruned UNet, and test the performance loss of the models after compression on the dataset Carvana [1].

  1. **UNet**: U-Net [18] is a deep convolutional neural network architecture that has proven highly effective for image segmentation tasks, especially in biomedical imaging. The network adopts an encoder-decoder structure. The encoder progressively downsamples feature maps to capture context, while the decoder upsamples these feature maps and concatenates them with corresponding feature maps from the encoder via skip connections to precisely localize objects.

  2. **Pruned UNet**: Specifically, the pruned UNet refers to the UNet model that has undergone Taylor-Rank Pruning [16] with FLOPs Regularization [11] (strength=0.001). This is essentially a channel pruning method. The bar chart below illustrates the percentage of channels pruned in each layer of the pruned UNet relative to the original UNet.



  3. **Dataset**: The Carvana Image Masking Challenge dataset [1] is a large-scale dataset specifically designed for image segmentation tasks, particularly in the automotive industry. It consists of a collection of high-resolution images of cars, along with corresponding pixel-level masks that delineate the precise boundaries of the vehicles. The Carvana dataset contains $5,088$ images with label masks in the training set and images without masks in the test set. Therefore, to evaluate model performance using the Dice metric, we randomly split $80\%$ of the dataset as the training set and $20\%$ as the validation set.

  4. **Evaluation**: The Dice coefficient, also known as the Dice similarity index, is a statistical measure used to gauge the similarity between two sets. It is particularly popular in the field of image segmentation to assess the accuracy of models by comparing the overlap between the predicted segmentation and the ground truth. Mathematically, the Dice coefficient is defined as

$$\texttt{Dice} = \frac{2|A \cap B|}{|A| + |B|}, \tag{12}$$

  where $A$ represents the predicted segmentation, and $B$ is the ground truth. The value of the Dice coefficient ranges from 0 to 1, where 1 indicates perfect agreement between two sets, and 0 indicates no overlap at all.

- **Details of experiments related to MobileNetV3**: We apply our compression method to both MobileNetV3 and the pruned MobileNetV3 optimized by automated gradual pruning algorithm [28], and evaluate the top-1 accuracy on the CIFAR-10 dataset [14].

    1. **MobileNetV3**: MobileNetV3 is a state-of-the-art convolutional neural network architecture designed specifically for mobile devices. Building upon the success of its predecessors, MobileNetV3 introduces novel design principles and optimization techniques to achieve a compelling balance between accuracy and efficiency.

    2. **Pruned MobileNetV3**: We employ the Automated Gradual Pruner (AGP) [28] technique to prune MobileNetV3, ultimately reducing the model size by approximately 2.68 times. The central idea of Automated Gradual Pruner (AGP) is to prune the model's weights gradually over time, allowing the network to adapt to the loss of parameters through continuous retraining. This technique is particularly effective in minimizing the performance degradation associated with aggressive pruning strategies.

    3. **Dataset**: The CIFAR-10 dataset [14] is a well-known benchmark in the field of machine learning and computer vision, consisting of $6,000$ $32 \times 32$ colored images belonging to 10 classes. The dataset is split into $50,000$ training images and $10,000$ testing images.

    4. **Evaluation**: The top-1 accuracy is an evaluation metric designed particularly for classification tasks, which primarily tests whether a model can make a correct single prediction for each input sample in a given task. Mathematically, let $\mathbf{x}_i$ represent an input instance, and let $\hat{y}_i = \arg\max_j \, p(y_j \mid \mathbf{x}_i)$ denote the predicted label, where $p(y_j \mid \mathbf{x}_i)$ is the predicted probability of the $j$-th class for the input $\mathbf{x}_i$. The true label is denoted by $y_i$. The top-1 accuracy is defined as the ratio of correctly predicted instances to the total number of instances, and can be expressed as the following formula:

    $$\text{Top-1 Accuracy} = \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(\hat{y}_i = y_i), \tag{13}$$

    where $N$ is the total number of instances, and $\mathbb{I}(\cdot)$ is the indicator function, which returns 1 if its argument is true and 0 otherwise.

- **Details of experiments related to LlaMA2-7B**: We apply our compression method to the LLaMA2-7B [21] and three other performant pruned variants of the LLaMA2-7B: Sheared-LLaMA-1.3B [25], TinyLLaMA [27], and LiteLLaMA. We then conduct performance tests on these models using various metrics, including SciQ [23], WinoGrande [19], ARC-E [5], 25-shot ARC-C [5], 10-shot Hellaswag [26], 32-shot BoolQ [4], 32-shot NQ [15], 5-shot MMLU [12], and Perplexity [24] to evaluate the effects of compression by using the package `lm-evaluation-harness` [9].

    1. **LlaMA2-7B [21]**: It is a state-of-the-art language model that has 7-billion parameters and was made publicly available by Meta AI. Employing the Transformer architecture, the model has been pre-trained on a massive text corpus, enabling it to exhibit superior performance on a multitude of natural language processing tasks, including text generation, question answering, machine translation, and so on.

    2. **Other comparative models**:
        - **Sheared-LlaMA-1.3B [25]**: It is a variant of LlaMA2-7B that uses a novel pruning algorithm: targeted structured pruning which can reduce a source model to a specified target architecture defined by the configurations of the existing pre-trained models. Additionally, this algorithm also introduces a dynamic batch loading strategy that proportionally loads training data from each domain based on its rate of loss reduction, thereby optimizing data utilization and facilitating overall performance improvements.
        - **TinyLlaMA [27]**: TinyLLaMA is a compact and efficient variant of the LLaMA2-7B model. It retains the same architecture and tokenizer as LLaMA2-7B but scales down to a Transformer decoder-only model with 1.1 billion parameters. This model is trained on up to 3 trillion tokens, with a specific focus on the behavior of smaller models when exposed to a substantially larger quantity of training tokens than those predicted by scaling laws.

– **LiteLlaMA**: LiteLlaMA is another variant of LlaMA2-7B with 460M parameters which is trained over 1T tokens on data from the RedPajama dataset [8].

3. **Evaluation**: In this study, most evaluation metrics for large language models (LLMs) are tested using the open-source package `lm-evaluation-harness` [9] which is a versatile and open-source Python package for benchmarking large language models (LLMs) across various tasks and datasets. The "(i)" after the name of an evaluation metric means "i-shot".

   – **SciQ** [23]: SciQ, a dataset of 13.7k multiple choice science exam questions, is a benchmark dataset specifically designed to evaluate a language model's ability to understand and reason about scientific questions. By presenting the model with carefully crafted questions that require a deep understanding of scientific concepts and principles, SciQ assesses a model's capacity to perform complex reasoning tasks.

   – **WinoGrande** [19]: Winogrande, a collection of 44k WSC-inspired problems, is a benchmark dataset specifically designed to evaluate a language model's ability to understand and reason about common sense knowledge. By presenting the model with carefully constructed sentences requiring the selection of the correct noun to complete a given context, Winogrande assesses a model's capacity to grasp nuanced semantic and pragmatic relationships.

   – **ARC-E & ARC-C (25)** [5]: ARC (AI2 Reasoning Challenge), including totally 7,787 natural, grade-school science questions, is a question set, text corpus, and baselines assembled to encourage AI research in advanced question answering, which is partitioned into a challenge set (ARC-C) and an easy set (ARC-E).

   – **HellaSwag (10)** [26]: HellaSwag, which stands for "Harder Endings, Longer contexts, and Low-shot Activities for Situations With Adversarial Generations," is designed to test commonsense natural language inference (NLI) in the context of physical situations.

   – **BoolQ (32)** [4]: BoolQ is a reading comprehension dataset of naturally occurring yes/no questions which are challenging and require a wide range of inference abilities to solve.

   – **NQ (32)** [15]: Natural Questions (NQ), including 307,372 training examples and 7,830 examples for development, contains real user questions issued to Google search, and answers found from Wikipedia by annotators, designed for the training and evaluation of automatic question answering systems.

   – **MMLU (5)** [12]: Massive Multitask Language Understanding (MMLU) benchmark is designed to evaluate a language model's ability to generalize across a wide range of tasks and domains. By assessing performance on a diverse set of 57 tasks, including elementary mathematics, US history, computer science, law and more, MMLU provides a comprehensive evaluation of a model's extensive world knowledge and problem solving ability.

   – **Perplexity** [24]: Perplexity (PPL) is used to measure the quality of a probability distribution or probabilistic model in predicting samples. A model with lower perplexity indicates better performance in predicting samples.

4. **More detailed information on all the models used for comparative testing**:

| Model | Parameter Precision | Number of Parameters | File Size |
|---|---|---|---|
| UNet [18] | FP32 | 13,395,329 | 52.4MB |
| Pruned-UNet | FP32 | 5,326,021 | 20.3MB |
| MobileNetV3 [13] | FP32 | 1,528,106 | 5.95MB |
| Pruned-MobileNetV3 | FP32 | 556,442 | 2.22MB |
| LlaMA2-7B [21] | FP16 | 6,738,415,616 | 12.50GB |
| Sheared-LlaMA-1.3B [25] | FP32 | 1,345,423,360 | 5.01GB |
| TinyLlaMA [27] | FP32 | 1,100,048,384 | 4.09GB |
| LiteLlaMA [1] | FP16 | 461,685,760 | 0.86GB |

---

[1] For more details, please visit the Hugging Face website of LiteLlaMA: `https://huggingface.co/ahxt/LiteLlama-460M-1T`.

# III. Supplementary Experimental Results

In this section, we present supplementary experimental results that are not put into the main body due to the limit of pages.

## 5.1 Compression Performance

Table 4: Comparison of Our Methods with Other Compression Variant Models of LlaMA2-7B on Individual Evaluation Dataset

| Model | File Size (GB) | SciQ (%) | WinoGrande (%) | ARC-E (%) | ARC- |
|---|---|---|---|---|---|
| LlaMA2-7B [21] | 12.50 | 94.00 | 68.98 | 76.30 | |
| Sheared-LlaMA-1.3B [25] | 5.01 (2.50×) | 87.30 | 58.09 | 60.98 | |
| TinyLlaMA [27] | 4.09 (3.06×) | 89.30 | 59.43 | 61.66 | |
| LiteLlaMA [1] | 0.86 (14.53×) | 75.10 | 52.64 | 47.85 | |
| LlaMA2-7B + HF | 4.80 (2.60×) | 93.70 | 69.77 | 75.59 | |
| Sheared-LlaMA-1.3B + HF | 0.98 (12.76×) | 87.90 | 58.88 | 60.48 | |
| TinyLlaMA + HF | 0.78 (16.03 ×) | 89.50 | 58.96 | 61.11 | |
| LiteLlaMA + HF | 0.39 (32.05×) | 73.00 | 54.22 | 44.78 | |

| Model | HellaSwag (10) (%) | BoolQ (32) (%) | NQ (32) (%) | MMLU (5) (%) | Ave |
|---|---|---|---|---|---|
| LlaMA2-7B [21] | 78.94 | 81.90 | 28.67 | 45.86 | |
| Sheared-LlaMA-1.3B [25] | 61.02 | 65.54 | 9.89 | 25.59 | |
| TinyLlaMA [27] | 62.48 | 62.91 | 12.52 | 26.79 | |
| LiteLlaMA | 38.41 | 57.09 | 1.77 | 26.11 | |
| LlaMA2-7B + HF | 77.17 | 80.92 | 25.65 | 43.04 | |
| Sheared-LlaMA-1.3B + HF | 60.56 | 63.76 | 8.98 | 24.68 | |
| TinyLlaMA + HF | 61.92 | 58.41 | 11.58 | 27.28 | |
| LiteLlaMA + HF | 37.56 | 56.57 | 1.16 | 26.61 | |

Table 5: The comparison of Perplexity (PPL) on the dataset *wikitext-2-raw-v1*

| Model | PPL |
|---|---|
| LlaMA2-7B | 5.19 |
| LlaMA2-7B + HF | 5.48 (+0.29) |
| Sheared-LlaMA-1.3B | 7.76 |
| Sheared-LlaMA-1.3B + HF | 7.99 (+0.23) |
| TinyLlaMA | 7.32 |
| TinyLlaMA + HF | 7.54 (+0.22) |
| LiteLlaMA | 20.22 |
| LiteLlaMA + HF | 23.89 (+3.67) |

Here, we show the compression performance of LLaMA2-7B and seven different variants based on LlaMA2-7B. First, in Table 4, we test eight important downstream tasks. By directly compressing LLaMA2-7B using our method (LLaMA2-7B + HF), a compression ratio of 2.6× is achieved with only a 0.99% decrease in the average score. When our method is combined with other LLaMA2-7B variants (Sheared-LlaMA-1.3B, TinyLlaMA, and LiteLlaMA) respectively, higher compression ratios can be achieved, with the change in the average score remaining within 1%, compared to the pruned ones. Moreover, Perplexity (PPL) is an important metric for evaluating the overall performance of large language models (LLMs). We conduct tests on the publicly available dataset wikitext-2-raw-v1

---

[1] For more details, please visit the Hugging Face website of LiteLlaMA: `https://huggingface.co/ahxt/LiteLlama-460M-1T`.

**My name is Clara and I am**

Prompt

Response

**LlaMA2-7B :**
My name is Clara and I am 13 years old. I live in the town of Boca Raton, Florida. I am a student at Spanish River High School.

**LlaMA2-7B + ET :**
My name is Clara and I am a 2015 graduate of the University of Southern California with a degree in Interactive Media and Game Design.

**Sheared-LlaMA-1.3B-Pruned :**
My name is Clara and I am a 20 year old student from the UK.

**Sheared-LlaMA-1.3B-Pruned + ET :**
My name is Clara and I am 15 years old. I am a student at the University of Applied Sciences in Vienna and I am currently studying to become a teacher.

**TinyLlaMA :**
My name is Clara and I am Clara.

**TinyLlaMA + ET :**
My name is Clara and I am to

**LiteLlaMA :**
My name is Clara and I am a teacher. I am a teacher of English and I am a teacher of Spanish.

**LiteLlaMA + ET :**
My name is Clara and I am a very happy and very intelligent girl.

**My name is Thomas and my main**

Prompt

Response

**LlaMA2-7B :**
My name is Thomas and my main interest is in creating and developing the best possible user experience for all users. I have a wide range of skills and experience. I have worked with everything from small businesses to large corporations.

**LlaMA2-7B + ET :**
My name is Thomas and my main goal is to be a professional soccer player.

**Sheared-LlaMA-1.3B-Pruned :**
My name is Thomas and my main hobby is playing the guitar. I have been playing for 10 years and I have been playing in bands for 5 years.

**Sheared-LlaMA-1.3B-Pruned + ET :**
My name is Thomas and my main hobby is playing video games. I have been playing video games since I was 5 years old and I have been playing them for 15 years.

**TinyLlaMA :**
My name is Thomas and my main name. I'm a 20 year old student who is currently studying for a BA in English Literature and Creative Writing.

**TinyLlaMA + ET :**
My name is Thomas and my main to

**LiteLlaMA :**
My name is Thomas and my main hobby is photography. I have been a professional photographer for over 20 years and have been fortunate to have worked with some of the best photographers in the world.

**LiteLlaMA + ET :**
My name is Thomas and my main hobby is photography. I have been a photographer for over 20 years and have been a member of the Photographic Society for over 10 years.

Figure 11: An example output from LlaMA2-7B, Sheared-LlaMA-1.3B, TinyLlaMA, LiteLlaMA and their compressed model by using our method. These results demonstrate that our compression method effectively preserves the models' ability to produce meaningful and grammatically correct text.

[2]. As shown in Table 5, the PPL of LLaMA2-7B+ET increases by only 0.29 compared to the original LLaMA2-7B model, whereas the PPL increases caused by the other three LLaMA2-7B variants (Sheared-LlaMA-1.3B, TinyLlaMA, LiteLlaMA) are 2.57, 2.13, and 15.03, respectively. This demonstrates that our method not only enhances the compression ratio but also effectively maintains the model's overall performance. Figure 11 presents two examples comparing the text outputs generated by eight different LLM models for a given prompt. The experimental data demonstrate that our method achieves an preferable compression ratio while having minimal impact on the model's performance across various downstream tasks. Moreover, it can also integrate effectively with other compression methods.

## 5.2 Compression Time

As shown in Table 6, we record the compression times for eight models using our method. The time to compress smaller models like UNet and MobileNetV3 is only 35 seconds, with Pruned-MobileNetV3 requiring only 8.83 seconds. The compression time for large language models is longer. For example, compressing LLaMA2-7B takes 53 minutes, while LiteLLaMA is compressed in just 6 minutes. The smaller the model's size, the faster the compression is.

Table 6: The compression time of our method on UNet, MobileNetV3 and LlaMA2-7B

| Model | Compression Time (s) |
|---|---|
| UNet + HF [18] | 30.00 |
| Pruned-UNet + HF | 34.60 |
| MobileNetV3 + HF | 11.63 |
| Pruned-MobileNetV3 + HF | 8.83 |
| **Model** | **Compression Time (min)** |
| LlaMA2-7B + HF | 53 |
| Sheared-LlaMA-1.3B + HF | 7 |
| TinyLlaMA + HF | 11 |
| LiteLlaMA + HF | 6 |

As shown in Figure 7, we conduct ablation experiments based on UNet and MobileNetV3 on the aforementioned two acceleration techniques to demonstrate their effectiveness in enhancing computational speed.

Table 7: The results of ablation study on acceleration techniques

| Model | K-D Tree | Matrix Operations | Compression Time (min) |
|---|---|---|---|
| | ✗ | ✗ | 231.6 |
| | ✓ | ✗ | 176.9 |
| UNet [18] | ✗ | ✓ | 54.8 |
| | ✓ | ✓ | **0.5** |
| | ✗ | ✗ | 52.29 |
| | ✓ | ✗ | 3.13 |
| MobileNetV3 [13] | ✗ | ✓ | 45.50 |
| | ✓ | ✓ | **0.19** |

Furthermore, we show that the concurrent application of k-d trees and matrix operations leads to a substantial enhancement in computational efficiency.

---

[2] For more details, please visit the website of dataset wikitext-2-raw-v1: `https://huggingface.co/datasets/Salesforce/wikitext`.

## 5.3 Inference Time

We test the inference time and the execution time of our method under different batch sizes on NVIDIA GeForce RTX 4060 and NVIDIA A6000 graphics cards, respectively.

Table 8: Timing results for different batch sizes on two devices using two methods with a UNet network.

| Batch Size | 4060 | | A6000 | |
|---|---|---|---|---|
| | Inference(s) | Our Method(s) | Inference(s) | Our Method(s) |
| 1 | 0.18 | 1.17 | 0.19 | 1.63 |
| 2 | 0.28 | 1.22 | 0.27 | 1.92 |
| 4 | 0.48 | 1.39 | 0.50 | 2.04 |
| 8 | 1.78 | 2.70 | 0.92 | 2.23 |
| 16 | 10.07 | 11.12 | 1.68 | 2.96 |
| 32 | / | / | 3.34 | 4.60 |

When performing the time test, the model is first initialized and loaded onto the GPU. Next, the start time is recorded, and the model evaluation function is executed 10 times in a loop. After each loop, the CUDA cache is manually cleared and Python's garbage collection mechanism is triggered to ensure that memory resource is fully released. After the loop ends, the end time is recorded, and the average time for each evaluation is calculated.

The Carvana dataset is used in the model evaluation function. This dataset is used for car image segmentation tasks and is widely used in image segmentation benchmarks. Each time the model is evaluated, only 1 batch of data is used for a test. We notice that when batch size increases, the inference time of our compressed model gets closer to that of the original model, which aligns with our expectations. This is because when batch size goes up, the inference time across layers is slow, which allows weights in later layers to be decompressed in time.

# V. Parameter Sensitivity

In this section, we use UNet as an example to conduct a parameter sensitivity test on three hyperparameters: $M, U$, and $l$ in our compression algorithm, where $M$ means the maximum value of categories that can be set for all layers of the model, $U$ means the list of the number of sample nodes in the center square, and $l$ means the length of the side of center square. As shown in Table 9, it can be observed that different values of $M, U$, and $l$ only moderately impact the model's performance and compression ratio, which means our algorithm is robust to hyperparameters.

Table 9: Parameter sensitivity results for different hyper-parameters on UNet.

| Model | Hyper-Parameters | | | Dice (%) | Compression Ratio |
|---|---|---|---|---|---|
| | Max Class | $U$ | $l$ | | |
| | - | - | - | 99.86 | - |
| | 1 | 225 | 0.02 | 99.81 | 7.10× |
| | 2 | 225 | 0.02 | 97.07 | 7.11× |
| | 3 | 225 | 0.02 | 90.11 | 7.11× |
| | 1 | 225 | 0.1 | 99.21 | 7.96× |
| | 2 | 225 | 0.1 | 96.46 | 7.93× |
| | **3** | **225** | **0.1** | **99.71** | **7.93×** |
| | 1 | 361 | 0.02 | 99.87 | 6.40× |
| | 2 | 361 | 0.02 | 97.90 | 6.41× |
| UNet | 3 | 361 | 0.02 | 99.39 | 6.40× |
| | 1 | 361 | 0.1 | 99.28 | 7.77× |
| | 2 | 361 | 0.1 | 99.10 | 7.75× |
| | 3 | 361 | 0.1 | 99.65 | 7.72× |
| | 1 | 225 / 361 | 0.02 | 99.89 | 7.11× |
| | 2 | 225 / 361 | 0.02 | 98.12 | 7.11× |
| | 3 | 225 / 361 | 0.02 | 99.58 | 7.11× |
| | 1 | 225 / 361 | 0.1 | 99.08 | 7.77× |
| | 2 | 225 / 361 | 0.1 | 98.85 | 7.75× |
| | 3 | 225 / 361 | 0.1 | 99.51 | 7.72× |

# References

[1] DanGill et al. Brian Shaler. Carvana image masking challenge, 2017.

[2] Simon Broadbent and J. M. Hammersley. Percolation processes. *Mathematical Proceedings of the Cambridge Philosophical Society*, 53:629 – 641, 1957.

[3] Sathya R Chitturi, Zhurun Ji, Alexander N Petsch, Cheng Peng, Zhantao Chen, Rajan Plumley, Mike Dunne, Sougata Mardanya, Sugata Chowdhury, Hongwei Chen, et al. Capturing dynamical correlations using implicit neural representations. *Nature Communications*, 14(1):5852, 2023.

[4] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint arXiv:1905.10044*, 2019.

[5] 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. *arXiv preprint arXiv:1803.05457*, 2018.

[6] Isaac P Cornfeld, Sergei Vasilevich Fomin, and Yakov Grigor'evic Sinai. *Ergodic theory*, volume 245. Springer Science & Business Media, 2012.

[7] Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, et al. Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature Machine Intelligence*, 5(3):220–235, 2023.

[8] Xue-Yong Fu, Md Tahmid Rahman Laskar, Elena Khasanova, Cheng Chen, and Shashi Bhushan TN. Tiny titans: Can smaller large language models punch above their weight in the real world for meeting summarization? *arXiv preprint arXiv:2402.00841*, 2024.

[9] L Gao, J Tow, B Abbasi, S Biderman, S Black, A DiPofi, C Foster, L Golding, J Hsu, A Le Noac'h, et al. A framework for few-shot language model evaluation. *URL https://zenodo. org/records/10256836*, 7, 2023.

[10] Geoffrey Grimmett. *Percolation*. Springer, 1991.

[11] Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks. In *Proceedings of the IEEE international conference on computer vision*, pages 1389–1397, 2017.

[12] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.

[13] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, Quoc V. Le, and Hartwig Adam. Searching for mobilenetv3, 2019.

[14] A Krizhevsky. Learning multiple layers of features from tiny images. *Master's thesis, University of Tront*, 2009.

[15] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466, 2019.

[16] Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning convolutional neural networks for resource efficient inference. *arXiv preprint arXiv:1611.06440*, 2016.

[17] Parikshit Ram and Kaushik Sinha. Revisiting kd-tree for nearest neighbor search. In *Proceedings of the 25th acm sigkdd international conference on knowledge discovery & data mining*, pages 1378–1388, 2019.

[18] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation, 2015.

[19] 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.

[20] Kenneth O Stanley, David B D'Ambrosio, and Jason Gauci. A hypercube-based encoding for evolving large-scale neural networks. *Artificial life*, 15(2):185–212, 2009.

[21] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023.

[22] M Mitchell Waldrop. The chips are down for moore's law. *Nature News*, 530(7589):144, 2016.

[23] Johannes Welbl, Nelson F Liu, and Matt Gardner. Crowdsourcing multiple choice science questions. *arXiv preprint arXiv:1707.06209*, 2017.

[24] Mengzhou Xia, Mikel Artetxe, Chunting Zhou, Xi Victoria Lin, Ramakanth Pasunuru, Danqi Chen, Luke Zettlemoyer, and Ves Stoyanov. Training trajectories of language models across scales. *arXiv preprint arXiv:2212.09803*, 2022.

[25] Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. Sheared llama: Accelerating language model pre-training via structured pruning, 2024.

[26] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.

[27] Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tinyllama: An open-source small language model, 2024.

[28] Michael Zhu and Suyog Gupta. To prune, or not to prune: exploring the efficacy of pruning for model compression. *arXiv preprint arXiv:1710.01878*, 2017.

[29] Xunyu Zhu, Jian Li, Yong Liu, Can Ma, and Weiping Wang. A survey on model compression for large language models. *arXiv preprint arXiv:2308.07633*, 2023.