

---

# Balancing Size and Sustainability: The Role of Compression in Large Language Models

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The rise of powerful Large Language Models like ChatGPT and GPT-4 has trans-  
2 formed AI across domains. However, their widespread use comes with significant  
3 environmental costs. To address this challenge, we propose a multi-pronged ap-  
4 proach, including LLM compression, resource optimization, and active monitoring.  
5 This paper focuses on evaluating compression methods for sustainable LLM de-  
6 ployment in enterprise settings.

## 7 1 Introduction

8 The emergence of powerful Large Language Models (LLMs), like ChatGPT and GPT-4 [10, 11],  
9 along with accessible and cost-effective APIs, has accelerated the adoption of Generative AI across  
10 diverse domains. Companies like Hugging Face have further democratized AI by making LLMs  
11 accessible to novice programmers and software developers. Beyond the consumer society, startups  
12 such as Glean, Scale, Yurts, OpenAI, and Cohere are at the vanguard of a revolution, propelling  
13 generative AI into enterprises through innovative on-premises deployments and personalized large  
14 language models. These LLMs, rooted in the Transformer architecture [14], are meticulously designed  
15 for language comprehension and boast an extensive output vocabulary comprising English words.  
16 Their prowess stems from training on colossal text corpora, approaching a staggering 2 trillion words,  
17 endowing them with the remarkable ability to generate coherent and meaningful text. The moniker  
18 "Large Language Models" is not just a label; it aptly encapsulates the scale of these models, with  
19 more modest variants like BERT [3] and Flan [1] encompassing between 100 million and 500 million  
20 parameters, while behemoths like GPT-3.5 and GPT-4 push the boundaries with an astonishing 150  
21 billion to 200 billion parameters.

22 While the widespread adoption of Large Language Models (LLMs) across domains is undeniably  
23 transformative, it brings to the forefront a pressing issue—the environmental cost of this technology’s  
24 proliferation. Beyond the monetary expenses associated with deploying LLMs, there is a significant  
25 environmental burden to consider. Take, for example, the resources required to host a single model  
26 with 100B (billion) parameters (=400GB memory) for standard inference tasks. A 100B parameter  
27 model requires a minimum of 25 consumer GPUs, wherein each GPU has a capacity of 16GB. The  
28 number of models an organization needs to host depends on factors like the expected number of  
29 concurrent users and their tolerance for latency. For instance, accommodating up to 2 seconds of  
30 latency for 100 users would necessitate hosting one model per 100 users, translating to 10 instances  
31 of the same model for a generative AI experience serving 1000 concurrent users. With a single model  
32 requiring 25 GPUs, this equates to a staggering 250 GPUs running continuously. The environmental  
33 implications are substantial, as a single A4000 GPU running for a day emits approximately 1.24kg  
34 CO<sub>2</sub>-eq [6], equivalent to driving an internal combustion engine car for 5.1 km. Scaling this to the  
35 operation of 250 GPUs to serve 1000 users on the platform for a single day results in an expected  
36 emission of 3108kg CO<sub>2</sub>-eq! which is equivalent to a car journey spanning 12,500km, slightly

37 more than a round trip from California to New York. Balancing the benefits of LLM's with their  
38 environmental impact poses a critical challenge in the era of Generative-AI expansion.

39 In the pursuit of sustainable deployment of Large Language Models (LLMs) and generative AI appli-  
40 cations, we advocate for a comprehensive, multi-pronged approach. Our strategy encompasses several  
41 key steps aimed at mitigating environmental impact while maintaining functional efficiency. First and  
42 foremost, we propose the development of mathematical methods to compress models to their sparser,  
43 low floating-point precision and lower parameter count counterparts of dense LLM's, while ensuring  
44 that their task performance remain uncompromised. Secondly, we tackle the scheduling challenge  
45 of efficiently allocating multiple models across minimal GPUs, introducing innovative sharding  
46 strategies to optimize resource utilization. Thirdly, we advocate a departure from a single, large  
47 foundational model to a mixture of smaller experts, significantly reducing GPU demands. Fourthly,  
48 we recommend the on-demand utilization of GPUs for hosting purposes and the implementation  
49 of predictive algorithms to estimate user traffic, minimizing the total GPU on-time. Lastly, we  
50 emphasize the importance of continuous monitoring, advocating for active tracking of GPU usage  
51 and its associated environmental costs within teams hosting LLMs and generative AI applications.  
52 This paper will focus on evaluating different compression methods for sustainable deployment of  
53 large language models, particularly when deployed at scale in enterprise settings.

## 54 2 Compression methods

55 Today, one of the most widely used commercial models is Meta's Llama-2-chat [13], introduced  
56 earlier this year. This model is pretrained on a massive dataset of 2 trillion tokens, roughly equivalent  
57 to about 500 billion words from the internet. Post pretraining, the model is further fine-tuned  
58 on multiple instruction datasets to enable it to provide precise answers to user queries (or user-  
59 instructions). The Llama-2 model comes in three different sizes: 7 billion parameters, 13 billion  
60 parameters, and 70 billion parameters, each demanding 28GB, 52GB, and 280GB of memory for their  
61 functioning. For hosting a single replica of the 7B, 13B, and 70B models, a minimum of 2, 4, and 18  
62 consumer GPUs, (with each consumer GPU having upto 16GB memory) is required, respectively. In  
63 order to maintain low latency as the number of concurrent users increase, it would be essential to  
64 scale the number of model replicas. We find that the optimal scaling estimate would be to have a  
65 single replica of the model for every 100 concurrent users, ensuring that user latency remains minimal  
66 while minimizing the idle GPU time.

67 In this paper, we are specifically interested in assessing the role of model compression on the  
68 sustainable deployment of LLMs. We will focus on 3 aspects of compression.

- 69 • **Model quantization:** wherein the parameters of the model are stored in lower precision  
70 floating points (e.g. conversion from fp-32 to fp-16 or int-4).
- 71 • **Model sparsification:** wherein the non-zero parameters of the model are reduced, i.e. a  
72 large fraction of the models' parameters are set to zero.
- 73 • **Model distillation:** wherein a large model is used as a "teacher" to train a smaller "student"  
74 model for task-specific use-cases.

### 75 2.1 Model quantization

76 The standard deployment of any LLM is done in such a way that each parameter of the model is  
77 stored as a 32 bit floating point (fp-32), wherein the parameter value is stored using 8 exponent bits,  
78 23 bits dedicated to the fraction and 1 bit for the sign (8+23+1=32 bits) [5], requiring 4 bytes per  
79 parameter. The estimates of model sizes provided in the earlier sections are arrived at by evaluating  
80 the product of 4 bytes per parameter with number of parameters (say, 7 billion) in the model (= 28GB).  
81 Quantization of models has been adopted as an effective method to quickly reduce the model size,  
82 across the ML industry. The common quantization levels used are: (i) 16 bits per parameter (fp-16),  
83 (ii) 19 bits per parameter (tfloat32), (iii) 16 per parameter (bfloat16), and (iii) 4 bits per parameter  
84 (int-4). That is, by quantizing the model to fp-16, or bfloat-16 the model size halves instantly, while  
85 quantizing to int-4 reduces the model size by 8 times! So, the Llama-2-7B model would require  
86 only 3.5GB of GPU memory while using an int-4 quantization and 14GB after fp-16 quantization,  
87 effectively reducing the number of GPUs required for hosting the model from 2GPU's to 1 GPU!

88 Also, depending on the hardware being used, the power consumption of fp16 and tfloat32 consume  
89 5x lesser power than fp32, while bfloat-16 consumes 9x lesser power than fp-32!

90 Although quantization works well for some use-cases, in practice, it has been observed that the naive  
91 reduction of floating point precision for every parameter in the model (especially to int-4) may result  
92 in the steep decline of functional performance. New methods, like, Sparse Quantized Representations  
93 [2], Reparameterized Ternary network (RTN) [7], have been proposed to identify parameters in the  
94 model that are more robust to quantization, in order to not sacrifice the generative performance of the  
95 quantized LLM obtained.

## 96 2.2 Model sparsification

97 Model sparsification refers to reducing the number of non-zero parameters in the LLM. The rationale  
98 behind model sparsification is that the non-zero weights in the LLM require upto 4 bytes per parameter  
99 (fp-32), while the zeroed out weights can be expressed with a single bit of information. The ML  
100 community has developed a wide array of sparsification routines, for reducing the number of non-  
101 zero weights in a neural network (or LLM). Some of the well known methods are: (1) Lottery  
102 ticket hypothesis [4], (2) Iterative pruning and retraining of neural networks [9, 8], (3) traversing  
103 functionally invariant paths (FIPs) [12], to name a few. Each of these techniques discover sparser  
104 counterparts of dense networks while maintaining the network (or LLMs) functional performance.

105 We subjected BERT networks to traversal along FIPs [12] and discovered sparse BERTs that retain  
106 their high performance on general language understanding tasks [15] (e.g. QQP, QNLI), and also  
107 maintained their perplexity on the Wikipedia dataset, although being 50% sparser 1A, B!

108 We also calculated the memory footprint of a few popular open-source LLMs (like Llama-2, Falcon)  
109 post sparsification by converting the underlying regular matrix (with a large fraction of zeroed out  
110 parameters) that make up the LLM into a sparse architecture (into Pytorch COO or CSR formats).  
111 In 1C, we find that a single instance of dense Falcon-180B requires more than 40 consumer GPUs  
112 (16GB memory each), while its 60% sparser counterpart requires only half the number of GPUs  
113 (~20) for hosting a single instance. In Fig. 1D, we find that the CO2 equivalent for running a single  
114 instance of Llama-2-70B for a single day can be mitigated by sparsifying the model as well as by  
115 choosing to host these models in appropriate datacenters across the globe, as it was observed that  
116 hosting models on AWS datacenters in Asian countries produce a much larger CO2 equivalent when  
117 compared to the US west coast.

118 Although, sparse networks should require much lesser memory for their inference, we **do not** observe  
119 this in practice. We observe that sparse networks only have lesser storage memory requirements  
120 (while stored using sparse matrices on a hard-drive), but during inference, they tend to require a  
121 similar amount of online GPU memory while being deployed for inferences. At times, for matrices  
122 that aren't "sparse enough", the online GPU memory requirement of a sparse matrix is much larger  
123 than that of its denser counterpart.

124 To actually harness the reduced computational footprint afforded by a sparsified model, we would  
125 need to enhance our hardware and build out routines at the assembly level to enable computation with  
126 matrices of mixed precision types. That is, to enable a single matrix to have elements with different  
127 memory types (say, fp-16, int-4, binary). Without these updates, it would be difficult for us to take  
128 advantage of LLMs that have unstructured sparsity (i.e. the non-zero parameters aren't contiguous  
129 blocks within a matrix, or within the LLM). However, with our current hardware and assembly level  
130 implementations, LLMs with structured sparsity can still be taken advantage of for reducing the  
131 environmental burden of wide-scale deployment.

## 132 2.3 Model distillation

133 Compared to the above mentioned methods of compressing models, the method of model distillation  
134 is significantly different, as it uses the original larger language model as the "teacher" model for  
135 training a much smaller "student" LLM (with any network architecture), on a task-specific use case.  
136 In this scenario, a synthetic dataset of inputs and target outputs are generated from the "teacher"  
137 LLM, which is subsequently used for training the student LLM via the supervised learning paradigm.

138 Distillation is extremely effective for building task-specific (smaller) models (or specialized experts).  
139 In many enterprise deployment settings, we find that the natural language tasks are very pointed,

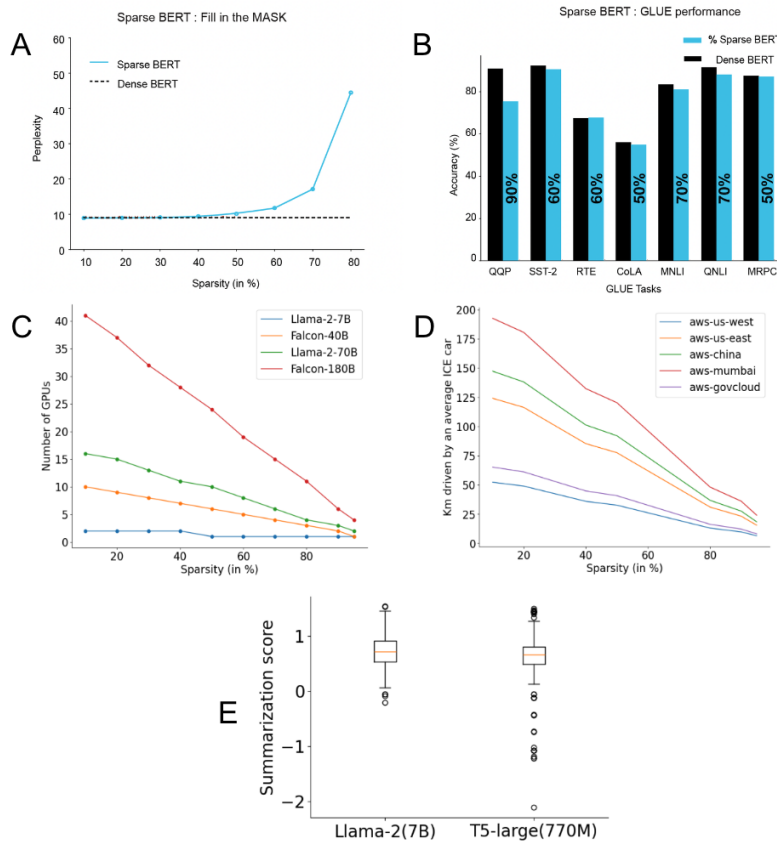


Figure 1: **Role of compression in sustainable deployment** (A) Perplexity score for sparse BERT (solid blue) and dense BERT (dashed blue). (B) GLUE performance of dense BERT (black) and p% sparse BERT (blue). (C) Quantifying memory requirement (in terms of GPU count) for a single replica of well known open-source models. (D) Quantifying CO2 emissions in terms of Km driven by an average ICE car [6] for hosting a single Llama-2-70B model for 24 hours. (E) Summarization score of the "teacher" Llama-2 (7B) and "student" T5-large (770M), post model distillation, on an evaluation dataset.

140 (for instance, summarization of an article), which makes distillation very effective. In Fig. 1E, we  
 141 observe that, distillation can be effectively used to shrink a Llama-2-7B (with 7B parameters) to a  
 142 T5-large (770M parameters) for a specific task like summarization without incurring a massive loss  
 143 in functional performance.

### 144 3 Discussion

145 In this paper, we've explored vital avenues for compressing large language models (LLMs) to enable  
 146 their sustainable deployment. Model quantization, sparsification, and distillation offer promising  
 147 ways to reduce memory and computational demands. Model quantization, achieved through lower-  
 148 precision formats, can significantly cut down model size and power consumption. However, a balance  
 149 must be struck to avoid sacrificing performance, necessitating advanced techniques. Sparsification of  
 150 LLMs can, theoretically speaking, lower GPU requirements and decrease CO2 emissions. However,  
 151 it currently doesn't reduce GPU memory usage during inference, highlighting the need for more  
 152 development on the hardware and assembly level routines. Model distillation, where smaller models  
 153 learn from larger ones, proves effective for task-specific applications, but aren't the best choice for  
 154 preserving the ability to perform a wide array of tasks. Choosing the right compression technique  
 155 depends on specific use cases and resource constraints. Future research and hardware improvements  
 156 will further enhance LLM deployment efficiency, ensuring their power is harnessed responsibly.

## References

- 157
- 158 [1] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li,  
159 Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned  
160 language models. *arXiv preprint arXiv:2210.11416*, 2022.
- 161 [2] Tim Dettmers, Ruslan Svirschevski, Vage Egiazarian, Denis Kuznedelev, Elias Frantar, Saleh  
162 Ashkboos, Alexander Borzunov, Torsten Hoefler, and Dan Alistarh. Spqr: A sparse-quantized  
163 representation for near-lossless llm weight compression. *arXiv preprint arXiv:2306.03078*,  
164 2023.
- 165 [3] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of  
166 deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*,  
167 2018.
- 168 [4] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable  
169 neural networks. *arXiv preprint arXiv:1803.03635*, 2018.
- 170 [5] William Kahan. Ieee standard 754 for binary floating-point arithmetic. *Lecture Notes on the*  
171 *Status of IEEE*, 754(94720-1776):11, 1996.
- 172 [6] Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. Quantifying  
173 the carbon emissions of machine learning. *arXiv preprint arXiv:1910.09700*, 2019.
- 174 [7] Yuhang Li, Xin Dong, Sai Qian Zhang, Haoli Bai, Yuanpeng Chen, and Wei Wang. Rtn: Repa-  
175 rameterized ternary network. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
176 volume 34, pages 4780–4787, 2020.
- 177 [8] Shiwei Liu, Decebal Constantin Mocanu, Amarsagar Reddy Ramapuram Matavalam, Yulong  
178 Pei, and Mykola Pechenizkiy. Sparse evolutionary deep learning with over one million artificial  
179 neurons on commodity hardware. *Neural Computing and Applications*, 33:2589–2604, 2021.
- 180 [9] Shiwei Liu, Lu Yin, Decebal Constantin Mocanu, and Mykola Pechenizkiy. Do we actually need  
181 dense over-parameterization? in-time over-parameterization in sparse training. In *International*  
182 *Conference on Machine Learning*, pages 6989–7000. PMLR, 2021.
- 183 [10] R OpenAI. Gpt-4 technical report. *arXiv*, pages 2303–08774, 2023.
- 184 [11] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin,  
185 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to  
186 follow instructions with human feedback. *Advances in Neural Information Processing Systems*,  
187 35:27730–27744, 2022.
- 188 [12] Guruprasad Raghavan and Matt Thomson. Engineering flexible machine learning systems by  
189 traversing functionally invariant paths in weight space. *arXiv preprint arXiv:2205.00334*, 2022.
- 190 [13] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,  
191 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open  
192 foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 193 [14] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
194 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*  
195 *processing systems*, 30, 2017.
- 196 [15] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman.  
197 Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv*  
198 *preprint arXiv:1804.07461*, 2018.