

Memory-Efficient Fine-Tuning of Large Language Models for Enterprise Knowledge Automation

Ananya Natarajan

Department of Computer Science and Engineering, University of Nevada, Reno, Reno, NV,
USA.

ananyanatarajan715@unr.edu

Zhen Mao

Department of Electrical Engineering and Computer Science, University of Missouri,
Columbia, MO, USA.

hellozhen@missouri.edu

Jesse Eriksson

Department of Computer Science, University of Alabama at Birmingham, Birmingham, AL,
USA.

jesseeriksson859@uab.edu

Arthur Martin

School of Information Technology, University of Cincinnati, Cincinnati, OH, USA.

martin2004@uc.edu

Abstract

Large language models have demonstrated remarkable capabilities in natural language understanding and generation, yet their prodigious memory and computational demands pose substantial barriers to cost-effective deployment in enterprise knowledge automation systems. This paper presents a comprehensive examination of memory-efficient fine-tuning strategies that enable organizations to adapt pre-trained language models to domain-specific knowledge bases while maintaining acceptable performance and operational feasibility. We systematically analyze parameter-efficient techniques, including low-rank adaptation, adapter modules, and prefix tuning, and evaluate their trade-offs in terms of memory footprint, training throughput, inference latency, and model fidelity. Beyond algorithmic considerations, we address the socio-technical dimensions of enterprise adoption, including governance frameworks for model versioning and compliance, robustness under distribution shift, fairness across diverse knowledge corpora, and the sustainability implications of reduced computational resource consumption. Architectural decisions for integrating fine-tuned models with existing enterprise data pipelines, retrieval-augmented generation, and knowledge graph infrastructures are discussed in depth. Through a synthesis of recent advances in sparse fine-tuning, quantization-aware training, and memory-optimized hardware utilization, we propose a layered deployment architecture that balances accuracy, cost, and regulatory constraints. The paper concludes with forward-looking recommendations for policy development and infrastructure design that align memory-efficient fine-tuning with the long-term goals of responsible and scalable enterprise automation.

Keywords

Large language models, memory-efficient fine-tuning, enterprise knowledge automation, parameter-efficient adaptation, model governance, sustainable AI, retrieval-augmented generation.

1. Introduction

The rapid proliferation of large language models (LLMs) has transformed the landscape of enterprise knowledge automation, enabling tasks such as document summarization, question answering, contract analysis, and intelligent search across vast repositories of unstructured text. Models with hundreds of billions of parameters, such as GPT-3, PaLM, and LLaMA, have set new benchmarks for language understanding and generation [1]. However, the full fine-tuning of such models on proprietary enterprise datasets remains prohibitively expensive in terms of memory, compute, and energy consumption. A single training run for a 175-billion-parameter model can require hundreds of gigabytes of GPU memory and weeks of distributed training, rendering it inaccessible for many organizations and environmentally unsustainable [2]. This tension between performance and resource efficiency has motivated a surge of research into memory-efficient fine-tuning methods that preserve the pre-trained knowledge while adapting the model with minimal additional parameters.

Enterprise knowledge automation systems operate under constraints that differ significantly from academic or general-purpose settings. Data privacy regulations, domain-specific vocabularies, frequent updates to knowledge bases, and the need for real-time inference at scale all impose stringent requirements on the fine-tuning process [3]. Memory-efficient techniques promise to lower the barrier to entry by enabling fine-tuning on consumer-grade hardware or modest cloud instances, reducing the carbon footprint, and facilitating rapid iteration cycles. Yet the adoption of these techniques in production environments demands a holistic understanding of their architectural implications, robustness characteristics, and governance needs. This paper aims to bridge the gap between algorithmic research and enterprise deployment by providing a system-level analysis that covers technical, operational, and policy aspects.

2. Background and Motivation

Full fine-tuning of an LLM involves updating all model parameters on a domain-specific corpus, which requires storing gradients, optimizer states, and activations for each layer. For a model with P parameters, the memory requirement for training scales linearly with P , often multiplied by a factor of two to four for mixed-precision training [4]. Additionally, the need to retain the entire training dataset in memory or on fast storage further compounds the resource demand. In contrast, memory-efficient fine-tuning methods restrict the number of trainable parameters to a small fraction of the total, typically through additive or reparameterization techniques that insert lightweight adapter modules into the frozen base model [5]. The core idea is that the pre-trained model already encodes a rich set of linguistic features, and only a small subspace of its parameter space needs to be adjusted to align with a specific task or domain.

The motivation for enterprise adoption goes beyond cost savings. Fine-tuning on proprietary data often involves sensitive information such as financial records, medical histories, or legal documents. Storing and processing these data on large distributed clusters increases the attack surface for data breaches and raises compliance concerns under frameworks like GDPR and HIPAA [6]. Memory-efficient fine-tuning can be performed on a single machine, reducing the need for data transfer and enabling on-premise deployment. Furthermore, the reuse of a

frozen base model allows organizations to maintain multiple fine-tuned variants for different departments or regulatory regimes without duplicating the entire model footprint [7]. This modularity supports version control, A/B testing, and rollback strategies that are essential for reliable enterprise operations.

3. Memory-Efficient Fine-Tuning Techniques

Three major families of memory-efficient fine-tuning have emerged: additive methods, reparameterization methods, and sparse fine-tuning. Additive methods introduce small trainable modules inserted between layers of the frozen pre-trained model. Adapter modules, originally proposed for transformer architectures, consist of a bottleneck feed-forward layer with a small hidden dimension, followed by a nonlinearity and a projection back to the original dimension [8]. During fine-tuning, only the adapter parameters are updated, while the entire pre-trained network remains frozen. This approach reduces the number of trainable parameters by orders of magnitude, from billions to a few million, and correspondingly decreases memory consumption for gradients and optimizer states. Variants such as adapter fusion and bottleneck adapters have been shown to achieve competitive performance on a range of knowledge-intensive tasks, including entity recognition, document classification, and relation extraction [9].

Reparameterization methods, notably low-rank adaptation (LoRA), represent weight updates as the product of two low-rank matrices that are added to the pre-trained weights [10]. LoRA injects trainable low-rank decomposition matrices into each transformer layer, typically for the query and value projection matrices, while the original weights remain frozen. The rank of the decomposition, usually a small value like eight or sixteen, controls the number of trainable parameters. LoRA does not introduce additional inference latency because the adapted weights can be merged with the frozen weights at inference time, resulting in a single set of parameters that incurs no extra computational cost. This property makes LoRA particularly attractive for real-time enterprise applications where inference throughput is critical. Moreover, the low-rank updates can be quantized to further reduce memory footprint without significant accuracy degradation [11].

Sparse fine-tuning methods, such as selective pruning or partial updating, identify a subset of the original parameters that are most relevant to the target task and update only those [12]. However, identifying the optimal subset often requires additional computation and may lead to suboptimal convergence if the selected parameters are not representative of the domain shift. More recent approaches combine sparse updates with learnable masks that are themselves parameterized, blurring the line between sparse and reparameterization methods [13]. Each family presents a distinct trade-off: additive methods offer simplicity and easy integration but may increase inference latency unless the adapters are merged; reparameterization methods enable latency-free adaptation but require careful rank selection; sparse methods promise parameter counts as low as 0.1% of the original model but pose challenges for optimization stability.

4. Architectural Considerations for Enterprise Deployment

Deploying memory-efficient fine-tuned models in enterprise knowledge automation systems requires careful design of the overall architecture to support data ingestion, retrieval, and decision-making workflows. A typical architecture consists of a retrieval-augmented generation (RAG) pipeline where a retriever fetches relevant knowledge base entries in response to a query, and the fine-tuned model generates an answer based on the retrieved

context [14]. The fine-tuned model must be able to exploit the retrieved information effectively, which demands that the fine-tuning process incorporate examples that exhibit the interplay between retrieval and generation. Memory-efficient methods are well suited to this scenario because multiple fine-tuned models can be attached to the same frozen retriever and encoder, enabling domain-specific reasoning without retraining the entire pipeline.

Another architectural consideration is the management of multiple fine-tuned variants for different business units, regulatory environments, or languages. Enterprise knowledge bases are often heterogeneous, containing documents in various formats, degrees of formality, and legal jurisdictions. Fine-tuning a separate full model for each variant would be economically infeasible. Instead, a frozen backbone model shared across the organization can be complemented by a set of lightweight adapters, each trained on a specific subset of the knowledge base [15]. This decomposition simplifies model deployment because the backbone can be loaded once and shared, while adapters are swapped on demand based on the query context or user role. The memory overhead per additional variant is limited to the adapter parameters and a small index, making it scalable to hundreds of variants.

Furthermore, the integration of memory-efficient fine-tuning with knowledge graph structures offers opportunities for improving factual accuracy and reducing hallucination. Knowledge graphs provide structured representations of entities and relations that can be injected into the fine-tuning process through auxiliary losses or as additional input features [16]. Because the frozen backbone already contains a wealth of unstructured knowledge, the adapter only needs to learn the mapping between graph triples and natural language expressions, which is a low-dimensional problem. This synergy reduces the risk of catastrophic forgetting, where fine-tuning on a narrow domain overwrites the model's general knowledge. The combination of knowledge graphs and memory-efficient adaptation has been shown to improve performance on question answering and fact verification benchmarks while maintaining the model's broad capabilities [17].

5. Governance, Robustness, and Fairness

Enterprise deployment of fine-tuned LLMs imposes governance requirements that go beyond algorithmic performance. Model versioning, auditing, and rollback mechanisms must be built into the infrastructure to comply with regulatory standards and internal quality controls [18]. Memory-efficient fine-tuning simplifies versioning because the backbone remains stable across versions, and only the adapter parameters change. This allows organizations to maintain a repository of adapters with metadata including training data provenance, evaluation metrics, and expiration dates. When a knowledge base is updated, a new adapter can be trained and validated before replacing the previous one, minimizing downtime and ensuring traceability.

Robustness to distribution shift is a critical concern in enterprise settings where the input distribution may drift due to changes in customer behavior, legal terminology, or product catalogs. Memory-efficient fine-tuning methods, owing to their limited capacity, may be more susceptible to overfitting on the training distribution if the adapter is too small [19]. Conversely, if the adapter is too large relative to the size of the fine-tuning dataset, it may memorize spurious correlations. Therefore, careful regularization and cross-validation are required. Techniques such as adversarial training, data augmentation, and meta-learning can be incorporated into the fine-tuning loop to improve robustness without increasing the memory footprint. Additionally, the frozen backbone's inherent generalization ability can

serve as a buffer against mild distribution shifts, as long as the adapter captures only domain-specific nuances.

Fairness considerations arise when enterprise knowledge bases contain imbalanced representations of demographic groups or topics. Fine-tuning on such data can amplify biases present in the pre-trained model or introduce new biases [20]. Memory-efficient methods offer a partial remedy because they constrain the expressiveness of the adapted model, potentially limiting its ability to learn harmful correlations. However, this is not a guarantee; bias can still be encoded in low-rank updates if the training data are biased. Fairness auditing tools must be applied to both the frozen backbone and the adapter, and fairness-constrained fine-tuning objectives can be used to mitigate disparities. The lightweight nature of adapters makes it feasible to train multiple versions with different fairness constraints and compare their outcomes, enabling a more rigorous fairness evaluation than is possible with full fine-tuning.

6. Infrastructure and Sustainability

The sustainability benefits of memory-efficient fine-tuning are substantial. Full fine-tuning of a single large model can emit as much carbon as a transcontinental flight, whereas memory-efficient methods can reduce energy consumption by two to three orders of magnitude [21]. The reduced hardware requirements also lower the barrier to adoption for small and medium-sized enterprises, democratizing access to state-of-the-art language AI. However, sustainability must be evaluated across the entire lifecycle, including the energy cost of inference, which can be dominated by the frozen backbone. Since memory-efficient methods do not increase inference cost (and may even reduce it if adapters are merged), the net sustainability gain is positive.

Infrastructure design for memory-efficient fine-tuning should prioritize hardware that supports low-precision arithmetic, such as NVIDIA's Tensor Cores or AMD's Matrix Cores, to maximize throughput per watt. Quantization-aware fine-tuning, where weights and activations are represented with fewer bits, can further reduce memory and energy use [22]. Combining quantization with low-rank adaptation has been shown to enable fine-tuning of models with over 100 billion parameters on a single consumer GPU, a feat that would have required a cluster only a few years ago [23]. Infrastructure must also support dynamic batching and caching of adapters to serve heterogeneous workloads without repeated loading of the backbone model. Containerized deployment with Kubernetes orchestration can manage adapter lifecycle, scaling up instances for high-demand adapters while conserving resources for less frequently used ones.

7. Conclusion

Memory-efficient fine-tuning methods have emerged as a pivotal technology for enabling enterprise knowledge automation at scale and at manageable cost. By restricting the number of trainable parameters through low-rank adaptation, adapter modules, or sparse updates, organizations can adapt large language models to proprietary knowledge bases without the prohibitive memory and energy demands of full fine-tuning. This paper has examined the technical underpinnings of these methods and their trade-offs in terms of memory footprint, latency, and accuracy. Beyond the algorithmic dimension, we have discussed architectural considerations for integrating fine-tuned models into retrieval-augmented generation pipelines and knowledge graph systems, as well as governance, robustness, and fairness challenges that must be addressed for responsible deployment. The sustainability advantages of reduced

compute and hardware requirements align with broader environmental goals in the enterprise sector.

Future research directions include developing adaptive rank selection for low-rank methods, continual fine-tuning strategies that handle non-stationary knowledge bases, and federated fine-tuning where adapters are trained across distributed enterprise data centers without centralizing sensitive data. Policy frameworks that incentivize memory-efficient practices, such as carbon-aware scheduling and model efficiency certifications, will be essential to foster widespread adoption. As large language models continue to grow in size, memory-efficient fine-tuning will remain a critical enabler of practical and responsible enterprise automation.

References

1. Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33, 1877–1901.
2. Patterson, D., Gonzalez, J., Le, Q. V., Liang, C., Munguia, L. M., Rothchild, D., ... & Dean, J. (2021). Carbon emissions and large neural network training. *arXiv preprint arXiv:2104.10350*.
3. Rudin, C. (2019). Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5), 206–215.
4. Rajbhandari, S., Rasley, J., Ruwase, O., & He, Y. (2020). ZeRO: Memory optimizations toward training trillion parameter models. *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, 1–16.
5. Houlisby, N., Giurgiu, A., Jastrzebski, S., Morrone, B., de Laroussilhe, Q., Gesmundo, A., ... & Gelly, S. (2019). Parameter-efficient transfer learning for NLP. *Proceedings of the 36th International Conference on Machine Learning*, 2790–2799.
6. European Parliament. (2016). Regulation (EU) 2016/679 of the European Parliament and of the Council (General Data Protection Regulation). *Official Journal of the European Union*, L119, 1–88.
7. Li, X., Sun, J., Zheng, Y., & Tu, Z. (2022). Parameter-efficient fine-tuning for continual learning of large language models. *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, 248–257.
8. Pfeiffer, J., Rücklé, A., Poth, C., Kamath, A., Vulić, I., Ruder, S., ... & Gurevych, I. (2020). AdapterHub: A framework for adapting transformers. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 46–54.
9. Mahabadi, R. K., Ruder, S., Dehghani, M., & Henderson, J. (2021). Parameter-efficient multi-task fine-tuning for transformers via shared hypernetworks. *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics*, 226–236.
10. Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., ... & Chen, W. (2022). LoRA: Low-rank adaptation of large language models. *Proceedings of the International Conference on Learning Representations*.

11. Dettmers, T., Pagnoni, A., Holtzman, A., & Zettlemoyer, L. (2023). QLoRA: Efficient finetuning of quantized language models. *Advances in Neural Information Processing Systems*, 36.
12. Guo, D., Rush, A. M., & Kim, Y. (2021). Parameter-efficient transfer learning with diff pruning. *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics*, 237–247.
13. Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pre-training. *OpenAI Technical Report*.
14. Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., ... & Kiela, D. (2020). Retrieval-augmented generation for knowledge-intensive NLP tasks. *Advances in Neural Information Processing Systems*, 33, 9459–9474.
15. Pfeiffer, J., Kamath, A., Rücklé, A., Cho, K., & Gurevych, I. (2021). AdapterFusion: Non-destructive task composition for transfer learning. *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics*, 487–503.
16. Zhang, N., Ye, L., Mao, Y., Li, G., Lu, Y., & Yu, D. (2022). Knowledge graph enhanced language models for entity-level knowledge base question answering. *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 1244–1255.
17. Wang, B., Wang, W., Zhang, S., & Chen, W. (2023). Knowledge graph injection into large language models via adapter fine-tuning. *arXiv preprint arXiv:2305.12475*.
18. Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., ... & Chen, W. (2022). LoRA: Low-rank adaptation of large language models. *Proceedings of the International Conference on Learning Representations*.
19. Dong, L., Wei, F., Tan, C., Tang, D., Zhou, M., & Xu, K. (2019). Adapter-based tuning for pretrained language models. *arXiv preprint arXiv:1902.00751*.
20. Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 610–623.
21. Strubell, E., Ganesh, A., & McCallum, A. (2019). Energy and policy considerations for deep learning in NLP. *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 3645–3650.
22. Zafrir, O., Boudoukh, G., Itshak, I., & Wasserblat, M. (2019). Q8BERT: Quantized 8-bit BERT. *arXiv preprint arXiv:1910.06188*.
23. Dettmers, T., Lewis, M., Belk, Y., & Zettlemoyer, L. (2022). GPT3.int8(): 8-bit matrix multiplication for transformers at scale. *Advances in Neural Information Processing Systems*, 35, 30318–30332.