

Review

Not peer-reviewed version

From Complexity to Simplicity: Advancements in Large Language Model Compression

Abdur Rashid Junaid^{*} and Idowu Callixtus^{*}

Posted Date: 25 December 2024

doi: 10.20944/preprints202412.2132.v1

Keywords: knowledge distillation; large language models; model compression; teacher-student models; neural networks; natural language processing; emergent behavior; efficient AI; multimodal distillation; adaptive distillation; AI scalability; task-specific fine-tuning; model optimization; sustainable AI



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Review

From Complexity to Simplicity: Advancements in Large Language Model Compression

Abdur Rashid Junaid ^{1,*} and Callixtus Idowu ^{2,*}

¹ Prince Mohammad Bin Fahd University, P.O. Box 1664 Al Khobar 31952 Kingdom of Saudi Arabia, Al-Khobar, Eastern Province, Saudi Arabia

² King Abdullah University of Science and Technology, Arabia, Saudi Arabia, 23955, Thuwa, Saudi Arabia

* Correspondence: abdur-rashid.junaid@pmu.edu.sa (A.R.J.); allixtus.idowu@kaust.edu.sa (C.I.)

Abstract: The rapid evolution of large language models (LLMs) has brought transformative advancements to natural language processing (NLP), enabling unprecedented performance in tasks such as machine translation, text generation, and conversational AI. However, the immense computational demands, memory requirements, and energy consumption of these models pose significant challenges for real-world deployment, particularly on resource-constrained devices and systems. Knowledge distillation has emerged as a pivotal technique for addressing these challenges, providing a framework for transferring knowledge from large, complex teacher models to smaller, efficient student models. This survey presents a comprehensive review of knowledge distillation methods tailored to LLMs, highlighting key advancements, applications, and unresolved challenges. We explore traditional distillation strategies, including logit matching and feature alignment, alongside contemporary approaches such as task-specific adaptations, attention map transfer, and progressive layer-by-layer distillation. Additionally, the integration of knowledge distillation with complementary compression techniques, such as quantization, pruning, and low-rank factorization, is examined, demonstrating their synergistic potential in optimizing LLMs for practical use. Applications of knowledge distillation span diverse domains, including edge computing, real-time systems, and fine-tuning for specialized tasks. The technique has facilitated the democratization of LLMs, enabling accessibility for organizations with limited computational resources. Despite these successes, challenges such as the efficient distillation of emergent behaviors, generalization in low-resource domains, and scalability for ultra-large teacher models remain significant barriers. Furthermore, the environmental impact and ethical considerations associated with LLM compression and deployment underscore the need for responsible innovation. Future directions in knowledge distillation research include the development of dynamic and adaptive distillation frameworks, automated processes leveraging neural architecture search, and benchmarks tailored to evaluate distillation outcomes comprehensively. By addressing these challenges, knowledge distillation can further enhance the efficiency, scalability, and inclusivity of LLMs, shaping the next generation of NLP systems.

Keywords: knowledge distillation; large language models; model compression; teacher-student models; neural networks; natural language processing; emergent behavior; efficient AI; multimodal distillation; adaptive distillation; AI scalability; task-specific fine-tuning; model optimization; sustainable AI

1. Introduction

Large language models (LLMs) have transformed the field of natural language processing (NLP), achieving state-of-the-art results in tasks such as machine translation, summarization, text classification, and conversational AI [1]. These models, characterized by their immense parameter counts and intricate architectures, have demonstrated remarkable capabilities, including in-context learning, zero-shot generalization, and high-quality text generation [2]. However, their exceptional performance comes at a cost: the training and deployment of LLMs demand substantial computational resources,

significant memory footprints, and considerable energy consumption [3,4], limiting their practical utility in resource-constrained environments and raising concerns about their environmental impact [5,6].

Knowledge distillation has emerged as a powerful technique to mitigate these challenges by transferring the knowledge embedded in a large, computationally intensive teacher model to a smaller, more efficient student model [7]. This process not only reduces the computational and storage requirements but also facilitates the deployment of LLMs on edge devices and in real-time applications without significant degradation in performance. In this survey, we provide an extensive review of the state-of-the-art in knowledge distillation for LLMs, covering the core methodologies, practical applications, and ongoing challenges in the field. We categorize and discuss various distillation approaches [8], including logit-based distillation, intermediate-layer alignment, and task-specific fine-tuning, while highlighting their strengths and limitations in different application scenarios[9].

The survey delves into real-world applications where knowledge distillation has enabled the efficient deployment of LLMs, such as low-latency conversational agents, personalized assistants, and mobile NLP applications [8,10]. Furthermore, we examine the integration of distillation with complementary model compression techniques, including pruning, quantization, and low-rank approximations, to maximize efficiency gains. We also explore emerging trends, such as multimodal distillation for models that process diverse data types and adaptive distillation frameworks tailored to specific tasks, datasets, or deployment constraints [11,12].

Despite the significant progress, knowledge distillation for LLMs faces several open challenges. These include the effective transfer of emergent behaviors such as reasoning and contextual understanding, scaling distillation techniques to handle ever-larger teacher models, and enhancing the generalization capabilities of student models across diverse domains and languages [13]. Additionally, the ethical and environmental implications of distillation warrant careful consideration, particularly in terms of bias transfer and the carbon footprint of training and distillation processes [14–16].

Looking ahead, we identify several promising directions for future research, including the development of automated and dynamic distillation pipelines, the exploration of multimodal and task-agnostic distillation strategies, and the design of sustainable frameworks to minimize the environmental impact of model compression. We also emphasize the need for standardized benchmarks and evaluation metrics to assess the effectiveness of distillation techniques comprehensively [17,18].

In conclusion, knowledge distillation represents a critical tool for addressing the scalability and efficiency challenges associated with LLMs. By enabling the creation of compact, high-performing models, it facilitates the democratization of cutting-edge NLP technologies, making them accessible to a broader audience. As the demand for scalable and environmentally sustainable AI systems continues to grow, knowledge distillation will undoubtedly play a pivotal role in shaping the next generation of LLMs and their applications.

Large Language Models (LLMs) have demonstrated remarkable capabilities across various natural language processing (NLP) tasks, ranging from machine translation and text summarization to sentiment analysis and question answering. These models, often consisting of billions of parameters, achieve state-of-the-art performance by leveraging vast datasets and sophisticated training techniques. However, the practical deployment of LLMs is fraught with challenges due to their significant computational and storage requirements. These challenges are particularly pronounced in edge devices and resource-constrained environments, where high latency and energy consumption become critical bottlenecks [19].

Model compression has emerged as a promising avenue to address these limitations. By reducing the size and complexity of LLMs, compression techniques aim to make these models more efficient without sacrificing their performance [20–22]. Among the various strategies for model compression, knowledge distillation stands out as an elegant and versatile approach. Knowledge distillation involves transferring the knowledge encoded in a large, pre-trained teacher model to a smaller, more efficient

student model. This process not only reduces the computational footprint of the model but also enables faster inference and deployment in real-world applications [23].

The core idea of knowledge distillation lies in leveraging the outputs or intermediate representations of the teacher model to guide the training of the student model [24]. These outputs may include the soft predictions of the teacher, which encapsulate rich information about the probability distribution over classes, or the hidden layer activations that capture intricate features of the input data. By aligning the student's learning process with the teacher's expertise, knowledge distillation ensures that the student model retains much of the teacher's performance despite its reduced size.

The application of knowledge distillation to LLMs, however, is not without its challenges. Unlike traditional machine learning models, LLMs exhibit unique characteristics, such as their ability to capture nuanced linguistic patterns and their reliance on extensive contextual information. Preserving these attributes during the distillation process requires innovative techniques and careful optimization [10]. Additionally, the scale of LLMs introduces computational hurdles in implementing effective distillation pipelines, necessitating advanced strategies to manage memory and processing resources.

Recent research in this domain has explored a wide array of approaches to enhance the effectiveness of knowledge distillation for LLM compression. These include task-specific distillation, where the student is optimized for a particular downstream task; multi-teacher distillation, which aggregates knowledge from multiple teacher models; and layer-wise distillation, which focuses on aligning specific layers of the student and teacher. Each of these methods addresses distinct aspects of the distillation challenge, contributing to a richer understanding of the field [25].

Moreover, innovative techniques such as self-distillation, where a model learns from its earlier checkpoints, and progressive distillation, which involves gradually reducing the teacher model's size, are gaining traction. These methods aim to enhance the transfer of knowledge by iteratively refining the student's capabilities. Additionally, leveraging attention mechanisms, contrastive learning, and reinforcement learning in the distillation process has shown promise in capturing the subtleties of LLMs [6].

Beyond the technical methodologies, evaluating the success of knowledge distillation requires robust metrics that account for both model performance and efficiency. Commonly used metrics include accuracy, latency, memory usage, and energy consumption, among others. These metrics provide a holistic view of the trade-offs involved in model compression and guide researchers in fine-tuning their approaches. Furthermore, the role of interpretability and explainability in compressed models is gaining attention, as understanding the behavior of smaller models becomes increasingly important in high-stakes applications.

The practical applications of knowledge distillation span various domains, including healthcare, finance, education, and more. For instance, deploying compressed LLMs in mobile devices enables personalized and privacy-preserving interactions, while resource-efficient models in healthcare can assist in diagnostics and decision-making [26,27]. These applications underscore the transformative potential of LLM compression in making advanced AI technologies more accessible and impactful [28,29].

This survey aims to provide a comprehensive overview of the advancements in knowledge distillation for LLM compression. We delve into the fundamental principles underlying this technique, explore diverse methodologies and their applications, and discuss the challenges and opportunities in the field. Furthermore, we examine the broader implications of LLM compression, including ethical considerations, environmental impact, and the democratization of AI technologies. By synthesizing insights from recent research, this survey seeks to equip researchers and practitioners with a deep understanding of the state-of-the-art in knowledge distillation and inspire future innovations in the field.

2. Related Work on Model Compression

Deep Neural Networks (DNNs) have revolutionized the field of artificial intelligence by achieving unprecedented accuracy across various tasks. However, the significant computational and storage requirements of these models present major challenges for their deployment in real-world scenarios, particularly on resource-constrained devices. Model compression techniques have emerged as a crucial area of research, aiming to address these limitations by reducing the size and complexity of DNNs while preserving their performance.

Several approaches have been proposed to compress DNNs [30], including parameter pruning, quantization, low-rank approximation, and knowledge distillation. Each of these techniques offers distinct advantages and is suited to different use cases:

1. **Parameter Pruning and Sparsity**: This technique involves identifying and removing redundant weights in a neural network. By enforcing sparsity, pruning reduces the number of parameters, leading to smaller model sizes and faster inference [31]. Methods such as structured pruning [32], which removes entire neurons or filters, and unstructured pruning, which removes individual weights, have shown significant promise [33].

2. **Quantization**: Quantization reduces the precision of model parameters and activations, often from 32-bit floating-point to lower bit-width representations such as 8-bit integers. This approach decreases the memory footprint and computational demands of the model, making it suitable for deployment on hardware accelerators like GPUs and TPUs [34].

3. **Low-Rank Approximation**: This method leverages the observation that weight matrices in neural networks often exhibit low-rank structures. By approximating these matrices with lower-rank representations, it is possible to achieve substantial reductions in model size without significant loss of accuracy [35].

4. **Knowledge Distillation**: Among the compression techniques, knowledge distillation occupies a unique position. Unlike other methods that directly modify the model architecture or parameters, knowledge distillation transfers the knowledge from a large, pre-trained teacher model to a smaller student model. The student is trained to mimic the outputs, intermediate representations, or attention patterns of the teacher, enabling it to achieve comparable performance with a fraction of the computational cost [36].

Knowledge distillation has gained widespread attention due to its versatility and effectiveness. It can be applied across various types of DNNs and tasks, including classification, object detection, and language modeling. Furthermore, it is often used in conjunction with other compression methods to achieve even greater efficiency [37]. For example, a pruned or quantized model can serve as the student in a distillation process, combining the benefits of multiple compression techniques.

The unique strength of knowledge distillation lies in its ability to encode rich, task-specific knowledge from the teacher into the student. This is particularly valuable in the context of LLMs, where the teacher model captures complex linguistic patterns and contextual information. By transferring this knowledge, distillation ensures that the student retains much of the teacher's functionality despite its reduced size [38].

In summary, model compression techniques, including pruning, quantization, low-rank approximation, and knowledge distillation, play a vital role in enabling the practical deployment of DNNs. Among these, knowledge distillation stands out for its ability to effectively bridge the gap between performance and efficiency, making it a cornerstone of modern compression strategies. The following sections delve deeper into the methodologies and advancements in knowledge distillation, highlighting its critical role in the compression of LLMs.

3. Techniques for Knowledge Distillation in LLMs

Knowledge distillation has proven to be a flexible and powerful approach for compressing Large Language Models (LLMs). Various techniques have been developed to effectively transfer knowledge from a large, pre-trained teacher model to a smaller, more efficient student model. This section

discusses key techniques that have emerged in recent years for distilling LLMs, focusing on their methodologies, applications, and effectiveness.

1. **Soft Label Distillation**: Soft label distillation is one of the most fundamental approaches in knowledge distillation. The teacher model generates probability distributions (soft labels) over the output classes, which encode richer information than hard labels alone [39]. These soft labels provide the student model with guidance on class relationships and decision boundaries. In the context of LLMs, soft label distillation can be applied to tasks such as text classification, where the probability distribution over classes captures nuanced semantic relationships [40].

2. **Feature-Based Distillation**: Feature-based distillation focuses on aligning the intermediate representations of the student model with those of the teacher. By transferring knowledge at the feature level, this method enables the student to capture detailed patterns and hierarchies learned by the teacher. Techniques such as layer-wise feature alignment and contrastive loss functions have been employed to enhance the effectiveness of feature-based distillation [41].

3. **Attention Transfer**: Attention mechanisms are a defining characteristic of LLMs, playing a crucial role in capturing contextual relationships in text. Attention transfer involves distilling the attention maps of the teacher into the student. By replicating the teacher's attention patterns, the student can better understand and process contextual dependencies. This technique has been shown to be particularly effective in transformer-based architectures [42].

4. **Task-Specific Distillation**: In task-specific distillation, the student model is optimized for a specific downstream task rather than general-purpose language understanding. The teacher provides task-specific guidance, such as task-specific logits or representations, which the student learns to mimic. This approach is especially useful when deploying LLMs for targeted applications such as sentiment analysis or machine translation [43].

5. **Multi-Teacher Distillation**: Multi-teacher distillation leverages the knowledge of multiple teacher models to guide the student. Each teacher may specialize in different aspects of the data or tasks, providing complementary knowledge to the student. Aggregating insights from multiple teachers can improve the robustness and generalization of the student model [44].

6. **Self-Distillation**: Self-distillation is a novel approach where a model learns from its own earlier checkpoints rather than a separate teacher model. This iterative process allows the model to refine its understanding and improve its performance over successive training iterations. Self-distillation is computationally efficient and aligns well with the progressive training of LLMs [22].

7. **Progressive Layer Distillation**: Progressive layer distillation involves transferring knowledge from the teacher to the student in a layer-by-layer manner. This gradual approach ensures that the student model can effectively learn from the teacher without being overwhelmed by the complexity of the full model. Progressive distillation is particularly beneficial for deep architectures, where transferring knowledge across all layers simultaneously may lead to suboptimal performance [45].

8. **Contrastive Learning for Distillation**: Contrastive learning techniques have been integrated into knowledge distillation to enhance the quality of knowledge transfer. By encouraging the student to learn representations that are similar to the teacher's for positive pairs and dissimilar for negative pairs, contrastive learning improves the alignment between the teacher and student representations. This approach is well-suited for tasks requiring fine-grained semantic understanding [46].

9. **Reinforcement Learning-Based Distillation**: Reinforcement learning (RL) has been used to optimize the distillation process, particularly for tasks involving sequential decision-making or structured outputs. RL-based distillation allows the student model to learn policies that mimic the teacher's behavior while optimizing for task-specific rewards [47].

Each of these techniques addresses different aspects of the distillation challenge, contributing to the versatility and effectiveness of knowledge distillation for LLM compression. By combining multiple techniques, researchers have achieved significant advances in reducing model size while maintaining or even enhancing performance on specific tasks. The next section explores the challenges

and open problems in knowledge distillation, highlighting areas for future research and innovation [48].

4. Future Works

Knowledge distillation for large language models (LLMs) is a rapidly evolving field with significant potential for innovation. While current research has demonstrated remarkable advancements, several open challenges and unexplored directions present opportunities for future work. Below, we outline key areas for further investigation and development.

4.1. Emergent Behavior Distillation

One of the defining features of large-scale language models is their ability to exhibit emergent behaviors, such as in-context learning, reasoning, and zero-shot generalization. These capabilities are often tied to the scale and architecture of the models, making them challenging to replicate in distilled student models. Future research should focus on designing distillation strategies that effectively transfer these emergent properties. This may include new loss functions, better alignment techniques for intermediate representations, or task-specific distillation frameworks that capture high-level reasoning.

4.2. Dynamic and Adaptive Distillation

Traditional distillation techniques rely on static processes, treating all layers or outputs of the teacher model equally. However, different tasks and datasets often require unique adaptations. Adaptive distillation, where the process dynamically adjusts based on the complexity of the task, the target application, or the computational budget, is an emerging area of interest. Research into reinforcement learning or meta-learning approaches for dynamic distillation pipelines could yield significant benefits[49].

4.3. Multimodal Distillation

With the rise of multimodal LLMs capable of processing text, images, and other data types simultaneously, extending knowledge distillation to these settings becomes increasingly important. Future work should address how to effectively distill multimodal knowledge, ensuring that student models retain the ability to integrate and reason across multiple data modalities while optimizing for size and efficiency[50].

4.4. Generalization Across Domains

LLMs are often required to perform across diverse domains and languages, but domain-specific generalization remains a challenge for student models. Investigating domain-aware distillation methods, including techniques that leverage cross-lingual or cross-domain pretraining, could enhance the robustness and versatility of distilled models. Additionally, techniques that minimize catastrophic forgetting during fine-tuning in new domains warrant further exploration [51,52].

4.5. Scalability of Distillation Techniques

As teacher models grow in scale, the computational cost of knowledge distillation increases significantly. Developing scalable distillation methods that can efficiently handle multi-billion-parameter teacher models is a critical direction for future research. Distributed and parallel distillation frameworks, as well as techniques that reduce memory overhead without sacrificing performance, could prove instrumental in this regard [9,53].

4.6. Integration with Other Compression Techniques

Knowledge distillation is often used in isolation, but its combination with other compression techniques such as pruning, quantization, and low-rank approximations has the potential to amplify

efficiency gains [54,55]. Future work should explore the interplay between these methods, aiming to develop unified frameworks that balance compression, accuracy, and computational feasibility [56].

Knowledge distillation (KD) and tensor decomposition are complementary techniques in the broader landscape of model compression, aimed at reducing the computational complexity and memory footprint of large language models (LLMs). KD focuses on transferring knowledge from a large teacher model to a smaller student model by aligning their outputs or internal representations, ensuring that the student model mimics the teacher's performance with reduced size and complexity. On the other hand, tensor decomposition techniques, such as matrix factorization, Tucker decomposition, or CP decomposition, compress models by factorizing weight tensors into lower-dimensional components, reducing the number of parameters and computations directly within the model architecture [57]. The relationship between the two approaches lies in their shared goal of efficiency: while KD optimizes the learning process to preserve task-specific performance, tensor decomposition structurally modifies the model to achieve efficiency. These techniques can be used synergistically, where a decomposed model serves as the student in KD, benefiting from both architectural compression and performance optimization. This combination has the potential to further reduce the computational cost of training and inference, making large models feasible for deployment in resource-constrained environments.

4.7. Environmental and Ethical Considerations

As LLMs become more prevalent, the environmental impact of training and distillation processes cannot be overlooked. Future research should aim to quantify and minimize the carbon footprint associated with model compression. Additionally, ethical considerations, such as bias transfer during distillation, must be addressed to ensure fairness and inclusivity in student models [58].

4.8. Automated Distillation Frameworks

Finally, the development of automated and end-to-end distillation frameworks represents a promising avenue. Such frameworks could leverage advances in neural architecture search (NAS) or hyperparameter optimization to automatically determine the best distillation strategy for a given task, dataset, or computational constraint. This would lower the barrier to entry for deploying efficient LLMs, making them more accessible to researchers and practitioners [59].

4.9. Evaluation Metrics and Benchmarks

Standardized benchmarks and evaluation metrics specific to knowledge distillation are still in their infancy. Future work should focus on creating comprehensive benchmarks that assess not only the accuracy and efficiency of distilled models but also their robustness, generalization, and ability to replicate emergent behaviors. Such benchmarks would provide clearer insights into the trade-offs involved in distillation and guide the development of more effective techniques [60].

In summary, the field of knowledge distillation for LLMs is ripe with opportunities for impactful research. By addressing these open challenges, future work can push the boundaries of what is achievable with compact and efficient models, enabling broader deployment and enhancing the accessibility of cutting-edge language technologies.

5. Conclusion

Knowledge distillation has become a pivotal strategy for addressing the computational and memory challenges posed by large language models (LLMs). It serves as a key enabler for deploying these models efficiently across diverse applications, including text generation, machine translation, and conversational AI. By transferring knowledge from large teacher models to smaller, resource-efficient student models, distillation facilitates reductions in model size and inference latency while preserving much of the teacher's performance benefits [61].

The field of knowledge distillation encompasses a wide array of methodologies, ranging from classic techniques such as logit matching to more advanced approaches involving task-specific adaptations, layer-wise alignment, and progressive strategies. These methods emphasize not only the replication of

teacher outputs but also the transfer of internal representations, capturing complex linguistic and semantic structures. This multifaceted approach underscores the effectiveness of distillation in enhancing the utility of LLMs [62].

Applications of knowledge distillation span domains requiring compact and efficient models, from real-time systems to resource-constrained edge deployments. By balancing computational efficiency with task performance, knowledge distillation enables broader accessibility of LLMs, democratizing their use among researchers and organizations with limited computational resources [63].

Despite significant progress, challenges persist. Distilling emergent properties like contextual reasoning and few-shot learning, which are closely tied to model scale, remains difficult. Furthermore, achieving robust generalization in low-resource and domain-specific settings is an open research area. As teacher models continue to grow in scale, the scalability of distillation methods becomes increasingly critical, necessitating innovative solutions to manage their computational complexity [6,8,53,64–68].

Future research directions promise to advance the field further. Combining distillation with other compression techniques such as quantization and pruning holds potential for synergistic gains. Adaptive and automated frameworks that tailor the distillation process to specific tasks, datasets, and resource constraints represent another exciting avenue. Additionally, exploring the ethical and environmental implications of model compression is essential to ensure the societal alignment of LLM technologies [21].

In summary, knowledge distillation stands as a cornerstone of modern NLP research, bridging the gap between the performance of large-scale models and the practical demands of deployment. By addressing both current challenges and emerging opportunities, distillation will continue to play a crucial role in shaping the evolution of efficient, scalable, and inclusive AI systems [47,69–73].

The rapid advancements in large language models (LLMs) have revolutionized natural language processing (NLP), enabling breakthroughs in diverse applications such as machine translation, text generation, and conversational AI. However, the increasing scale and complexity of these models pose significant challenges in terms of computational requirements, memory footprint, and energy consumption. Knowledge distillation has emerged as a critical technique for addressing these issues, offering a pathway to reduce model size and inference latency while retaining most of the performance benefits of the original models.

In this survey, we have provided a comprehensive overview of the state-of-the-art methodologies, applications, and challenges associated with knowledge distillation for LLMs. The surveyed techniques highlight the diversity and depth of this field, ranging from traditional approaches focusing on logit matching and feature-based distillation to more sophisticated methods leveraging task-specific adaptations, intermediate layer alignment, and progressive distillation strategies. Furthermore, recent research emphasizes the importance of aligning student models not only with the outputs of the teacher but also with its internal representations, fostering a deeper understanding of linguistic and semantic structures [61].

Applications of knowledge distillation in LLMs span a wide array of domains. These include creating compact models for deployment on edge devices, optimizing model efficiency for real-time systems, and enabling resource-efficient fine-tuning for specialized tasks. By achieving a balance between computational efficiency and task performance, knowledge distillation has facilitated the democratization of LLMs, making them accessible to researchers and organizations with limited computational resources [62].

Despite the impressive progress, significant challenges remain unresolved. One notable challenge is the distillation of emergent capabilities observed in the largest models, such as few-shot learning and contextual reasoning, which are often tied to the scale of the model. Additionally, the generalization capabilities of student models in low-resource and domain-specific settings require further investigation. Another critical issue is the scalability of distillation techniques as teacher models continue to grow in

size, necessitating innovative approaches that can handle the increasing computational complexity [63].

Future research directions in this field are promising and multifaceted. First, exploring the integration of knowledge distillation with other model compression techniques, such as quantization and pruning, may yield synergistic benefits. Second, developing automated and adaptive distillation frameworks that dynamically adjust the distillation process based on the task, dataset, and computational constraints could further enhance efficiency. Third, there is an opportunity to investigate the ethical and environmental implications of LLM deployment and compression, ensuring that the benefits of distillation are aligned with broader societal goals [64].

In conclusion, knowledge distillation has established itself as a cornerstone of research aimed at making LLMs more efficient and accessible. By reducing the barriers to deploying these powerful models in real-world applications, knowledge distillation not only addresses pressing computational challenges but also opens avenues for innovation and inclusivity in AI research and development [21]. As the demand for scalable AI systems grows, continued advancements in knowledge distillation will play a pivotal role in shaping the future of NLP and AI at large.

References

1. Agarwal, R.; Vieillard, N.; Zhou, Y.; Stanczyk, P.; Garea, S.R.; Geist, M.; Bachem, O. On-Policy Distillation of Language Models: Learning from Self-Generated Mistakes. In Proceedings of the The Twelfth International Conference on Learning Representations, 2024.
2. Burges, C.; Shaked, T.; Renshaw, E.; Lazier, A.; Deeds, M.; Hamilton, N.; Hullender, G. Learning to rank using gradient descent. In Proceedings of the Proceedings of the 22nd International Conference on Machine Learning, New York, NY, USA, 2005; ICML '05, p. 89–96. <https://doi.org/10.1145/1102351.1102363>.
3. Jiang, Y.; Chan, C.; Chen, M.; Wang, W. Lion: Adversarial Distillation of Closed-Source Large Language Model. *arXiv preprint arXiv:2305.12870* 2023.
4. Li, L.; Xie, Z.; Li, M.; Chen, S.; Wang, P.; Chen, L.; Yang, Y.; Wang, B.; Kong, L. Silkie: Preference Distillation for Large Visual Language Models. *arXiv preprint arXiv:2312.10665* 2023.
5. Allen-Zhu, Z.; Li, Y. Towards understanding ensemble, knowledge distillation and self-distillation in deep learning. *arXiv preprint arXiv:2012.09816* 2020.
6. Sachan, D.; Lewis, M.; Joshi, M.; Aghajanyan, A.; Yih, W.t.; Pineau, J.; Zettlemoyer, L. Improving Passage Retrieval with Zero-Shot Question Generation. In Proceedings of the Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing; Goldberg, Y.; Kozareva, Z.; Zhang, Y., Eds., Abu Dhabi, United Arab Emirates, 2022; pp. 3781–3797. <https://doi.org/10.18653/v1/2022.emnlp-main.249>.
7. Bai, Y.; Kadavath, S.; Kundu, S.; Askell, A.; Kernion, J.; Jones, A.; Chen, A.; Goldie, A.; Mirhoseini, A.; McKinnon, C.; et al. Constitutional AI: Harmlessness from AI Feedback, 2022, [arXiv:cs.CL/2212.08073].
8. Gangal, V.; Feng, S.Y.; Alikhani, M.; Mitamura, T.; Hovy, E. Nareor: The narrative reordering problem. In Proceedings of the Proceedings of the AAAI Conference on Artificial Intelligence, 2022, Vol. 36, pp. 10645–10653.
9. Wan, Z.; Wang, X.; Liu, C.; Alam, S.; Zheng, Y.; Liu, J.; Qu, Z.; Yan, S.; Zhu, Y.; Zhang, Q.; et al. Efficient Large Language Models: A Survey, 2024, [arXiv:cs.CL/2312.03863].
10. Lai, J.; Gan, W.; Wu, J.; Qi, Z.; Yu, P.S. Large Language Models in Law: A Survey. *arXiv preprint arXiv:2312.03718* 2023.
11. Dai, H.; Liu, Z.; Liao, W.; Huang, X.; Cao, Y.; Wu, Z.; Zhao, L.; Xu, S.; Liu, W.; Liu, N.; et al. AugGPT: Leveraging ChatGPT for Text Data Augmentation, 2023, [arXiv:cs.CL/2302.13007].
12. Sachan, D.S.; Lewis, M.; Yogatama, D.; Zettlemoyer, L.; Pineau, J.; Zaheer, M. Questions Are All You Need to Train a Dense Passage Retriever. *Transactions of the Association for Computational Linguistics* 2023, 11, 600–616.
13. Sanh, V.; Debut, L.; Chaumond, J.; Wolf, T. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108* 2019.
14. Liu, W.; Zeng, W.; He, K.; Jiang, Y.; He, J. What Makes Good Data for Alignment? A Comprehensive Study of Automatic Data Selection in Instruction Tuning, 2023, [arXiv:cs.CL/2312.15685].
15. Kiesel, J.; Alshomary, M.; Handke, N.; Cai, X.; Wachsmuth, H.; Stein, B. Identifying the Human Values behind Arguments. In Proceedings of the Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers); Muresan, S.; Nakov, P.; Villavicencio, A., Eds., Dublin, Ireland, 2022; pp. 4459–4471. <https://doi.org/10.18653/v1/2022.acl-long.306>.

16. Kang, M.; Lee, S.; Baek, J.; Kawaguchi, K.; Hwang, S.J. Knowledge-Augmented Reasoning Distillation for Small Language Models in Knowledge-Intensive Tasks, 2023, [arXiv:cs.CL/2305.18395].
17. Meng, R.; Liu, Y.; Yavuz, S.; Agarwal, D.; Tu, L.; Yu, N.; Zhang, J.; Bhat, M.; Zhou, Y. AugTrieve: Unsupervised Dense Retrieval by Scalable Data Augmentation, 2023, [arXiv:cs.CL/2212.08841].
18. Magister, L.C.; Mallinson, J.; Adamek, J.; Malmi, E.; Severyn, A. Teaching Small Language Models to Reason. In Proceedings of the Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers); Rogers, A.; Boyd-Graber, J.; Okazaki, N., Eds., Toronto, Canada, 2023; pp. 1773–1781. <https://doi.org/10.18653/v1/2023.acl-short.151>.
19. Gu, Y.; Dong, L.; Wei, F.; Huang, M. MiniLLM: Knowledge Distillation of Large Language Models. In Proceedings of the The Twelfth International Conference on Learning Representations, 2024.
20. Li, G.; Hammoud, H.A.A.K.; Itani, H.; Khizbullin, D.; Ghanem, B. Camel: Communicative agents for" mind" exploration of large scale language model society. *arXiv preprint arXiv:2303.17760* **2023**.
21. Lee, Y.S.; Sultan, M.; El-Kurdi, Y.; Naseem, T.; Munawar, A.; Florian, R.; Roukos, S.; Astudillo, R. Ensemble-Instruct: Instruction Tuning Data Generation with a Heterogeneous Mixture of LMs. In Proceedings of the Findings of the Association for Computational Linguistics: EMNLP 2023; Bouamor, H.; Pino, J.; Bali, K., Eds., Singapore, 2023; pp. 12561–12571. <https://doi.org/10.18653/v1/2023.findings-emnlp.836>.
22. Ye, S.; Jo, Y.; Kim, D.; Kim, S.; Hwang, H.; Seo, M. SelfFee: Iterative Self-Revising LLM Empowered by Self-Feedback Generation. Blog post, 2023.
23. Cao, H.; Liu, Z.; Lu, X.; Yao, Y.; Li, Y. InstructMol: Multi-Modal Integration for Building a Versatile and Reliable Molecular Assistant in Drug Discovery. *CoRR* **2023**, *abs/2311.16208*, [2311.16208]. <https://doi.org/10.48550/ARXIV.2311.16208>.
24. Summers, T.; Marino, K.; Ahuja, A.; Fergus, R.; Dasgupta, I. Distilling internet-scale vision-language models into embodied agents. In Proceedings of the Proceedings of the 40th International Conference on Machine Learning. JMLR.org, 2023, ICML/23.
25. Lou, R.; Zhang, K.; Xie, J.; Sun, Y.; Ahn, J.; Xu, H.; Su, Y.; Yin, W. MUFFIN: Curating Multi-Faceted Instructions for Improving Instruction-Following, 2023, [arXiv:cs.CL/2312.02436].
26. Liu, W.; Li, G.; Zhang, K.; Du, B.; Chen, Q.; Hu, X.; Xu, H.; Chen, J.; Wu, J. Mind's Mirror: Distilling Self-Evaluation Capability and Comprehensive Thinking from Large Language Models, 2023, [arXiv:cs.CL/2311.09214].
27. Liang, K.J.; Hao, W.; Shen, D.; Zhou, Y.; Chen, W.; Chen, C.; Carin, L. MixKD: Towards Efficient Distillation of Large-scale Language Models. In Proceedings of the 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021.
28. Liu, H.; Li, C.; Li, Y.; Lee, Y.J. Improved Baselines with Visual Instruction Tuning, 2023, [arXiv:cs.CV/2310.03744].
29. Kim, M.; Lee, S.; Lee, J.; Hong, S.; Chang, D.S.; Sung, W.; Choi, J. Token-Scaled Logit Distillation for Ternary Weight Generative Language Models. *arXiv preprint arXiv:2308.06744* **2023**.
30. Wei, J.; Tay, Y.; Bommasani, R.; Raffel, C.; Zoph, B.; Borgeaud, S.; Yogatama, D.; Bosma, M.; Zhou, D.; Metzler, D.; et al. Emergent Abilities of Large Language Models. *Trans. Mach. Learn. Res.* **2022**, 2022.
31. Li, M.; Chen, L.; Chen, J.; He, S.; Gu, J.; Zhou, T. Selective Reflection-Tuning: Student-Selected Data Recycling for LLM Instruction-Tuning. 2024.
32. Zniyed, Y.; Nguyen, T.P.; et al. Efficient tensor decomposition-based filter pruning. *Neural Networks* **2024**, *178*, 106393.
33. Li, Y.; Yu, Y.; Zhang, Q.; Liang, C.; He, P.; Chen, W.; Zhao, T. LoSparse: Structured Compression of Large Language Models based on Low-Rank and Sparse Approximation, 2023, [arXiv:cs.LG/2306.11222].
34. Kim, Y.J.; Henry, R.; Fahim, R.; Awadalla, H.H. FineQuant: Unlocking Efficiency with Fine-Grained Weight-Only Quantization for LLMs, 2023, [arXiv:cs.LG/2308.09723].
35. Radford, A.; Wu, J.; Child, R.; Luan, D.; Amodei, D.; Sutskever, I.; et al. Language models are unsupervised multitask learners. *OpenAI blog* **2019**, *1*, 9.
36. Sun, W.; Xie, R.; Zhang, J.; Zhao, W.X.; Lin, L.; Wen, J.R. Distillation is All You Need for Practically Using Different Pre-trained Recommendation Models. *arXiv preprint arXiv:2401.00797* **2024**.
37. Plummer, B.A.; Wang, L.; Cervantes, C.M.; Caicedo, J.C.; Hockenmaier, J.; Lazebnik, S. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In Proceedings of the Proceedings of the IEEE international conference on computer vision, 2015, pp. 2641–2649.

38. Schick, T.; Schütze, H. Generating Datasets with Pretrained Language Models. In Proceedings of the Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing; Moens, M.F.; Huang, X.; Specia, L.; Yih, S.W.t., Eds., Online and Punta Cana, Dominican Republic, 2021; pp. 6943–6951. <https://doi.org/10.18653/v1/2021.emnlp-main.555>.
39. Longpre, S.; Lu, Y.; Tu, Z.; DuBois, C. An Exploration of Data Augmentation and Sampling Techniques for Domain-Agnostic Question Answering. In Proceedings of the Proceedings of the 2nd Workshop on Machine Reading for Question Answering; Fisch, A.; Talmor, A.; Jia, R.; Seo, M.; Choi, E.; Chen, D., Eds., Hong Kong, China, 2019; pp. 220–227. <https://doi.org/10.18653/v1/D19-5829>.
40. Wang, Z.; Yu, A.W.; Firat, O.; Cao, Y. Towards Zero-Label Language Learning, 2021, [\[arXiv:cs.CL/2109.09193\]](https://arxiv.org/abs/2109.09193).
41. Chen, B.; Shu, C.; Shareghi, E.; Collier, N.; Narasimhan, K.; Yao, S. FireAct: Toward Language Agent Fine-tuning, 2023, [\[arXiv:cs.CL/2310.05915\]](https://arxiv.org/abs/2310.05915).
42. Sun, W.; Yan, L.; Ma, X.; Wang, S.; Ren, P.; Chen, Z.; Yin, D.; Ren, Z. Is ChatGPT Good at Search? Investigating Large Language Models as Re-Ranking Agents, 2023, [\[arXiv:cs.CL/2304.09542\]](https://arxiv.org/abs/2304.09542).
43. Liang, Y.; Wu, C.; Song, T.; Wu, W.; Xia, Y.; Liu, Y.; Ou, Y.; Lu, S.; Ji, L.; Mao, S.; et al. TaskMatrix.AI: Completing Tasks by Connecting Foundation Models with Millions of APIs, 2023, [\[arXiv:cs.AI/2303.16434\]](https://arxiv.org/abs/2303.16434).
44. Mukherjee, S.; Mitra, A.; Jawahar, G.; Agarwal, S.; Palangi, H.; Awadallah, A. Orca: Progressive learning from complex explanation traces of gpt-4. *arXiv preprint arXiv:2306.02707* **2023**.
45. Chen, Z.; Deng, Y.; Yuan, H.; Ji, K.; Gu, Q. Self-Play Fine-Tuning Converts Weak Language Models to Strong Language Models, 2024, [\[arXiv:cs.LG/2401.01335\]](https://arxiv.org/abs/2401.01335).
46. Wu, T.; Luo, L.; Li, Y.F.; Pan, S.; Vu, T.T.; Haffari, G. Continual Learning for Large Language Models: A Survey. *arXiv preprint arXiv:2402.01364* **2024**.
47. Ren, X.; Wei, W.; Xia, L.; Su, L.; Cheng, S.; Wang, J.; Yin, D.; Huang, C. Representation Learning with Large Language Models for Recommendation, 2023, [\[arXiv:cs.IR/2310.15950\]](https://arxiv.org/abs/2310.15950).
48. Nakano, R.; Hilton, J.; Balaji, S.; Wu, J.; Ouyang, L.; Kim, C.; Hesse, C.; Jain, S.; Kosaraju, V.; Saunders, W.; et al. WebGPT: Browser-assisted question-answering with human feedback, 2022, [\[arXiv:cs.CL/2112.09332\]](https://arxiv.org/abs/2112.09332).
49. Padmanabhan, S.; Onoe, Y.; Zhang, M.J.; Durrett, G.; Choi, E. Propagating Knowledge Updates to LMs Through Distillation. *arXiv preprint arXiv:2306.09306* **2023**.
50. Cai, Z.; Tao, C.; Shen, T.; Xu, C.; Geng, X.; Lin, X.A.; He, L.; Jiang, D. HypeR: Multitask Hyper-Prompted Training Enables Large-Scale Retrieval Generalization. In Proceedings of the The Eleventh International Conference on Learning Representations, 2022.
51. Qian, C.; Han, C.; Fung, Y.; Qin, Y.; Liu, Z.; Ji, H. CREATOR: Tool Creation for Disentangling Abstract and Concrete Reasoning of Large Language Models. In Proceedings of the Findings of the Association for Computational Linguistics: EMNLP 2023; Bouamor, H.; Pino, J.; Bali, K., Eds., Singapore, 2023; pp. 6922–6939. <https://doi.org/10.18653/v1/2023.findings-emnlp.462>.
52. Qiu, L.; Zhao, Y.; Li, J.; Lu, P.; Peng, B.; Gao, J.; Zhu, S.C. Valuenet: A new dataset for human value driven dialogue system. In Proceedings of the Proceedings of the AAAI Conference on Artificial Intelligence, 2022, Vol. 36, pp. 11183–11191.
53. Wu, C.; Zhang, X.; Zhang, Y.; Wang, Y.; Xie, W. PMC-LLaMA: Further Finetuning LLaMA on Medical Papers. *CoRR* **2023**, *abs/2304.14454*, [\[2304.14454\]](https://arxiv.org/abs/2304.14454). <https://doi.org/10.48550/ARXIV.2304.14454>.
54. Lu, D.; Wu, H.; Liang, J.; Xu, Y.; He, Q.; Geng, Y.; Han, M.; Xin, Y.; Xiao, Y. BBT-Fin: Comprehensive Construction of Chinese Financial Domain Pre-trained Language Model, Corpus and Benchmark. *CoRR* **2023**, *abs/2302.09432*, [\[2302.09432\]](https://arxiv.org/abs/2302.09432). <https://doi.org/10.48550/ARXIV.2302.09432>.
55. Yu, F.; Gao, A.; Wang, B. Outcome-supervised Verifiers for Planning in Mathematical Reasoning. *CoRR* **2023**, *abs/2311.09724*, [\[2311.09724\]](https://arxiv.org/abs/2311.09724). <https://doi.org/10.48550/ARXIV.2311.09724>.
56. Feng, S.Y.; Gangal, V.; Wei, J.; Chandar, S.; Vosoughi, S.; Mitamura, T.; Hovy, E. A survey of data augmentation approaches for NLP. *arXiv preprint arXiv:2105.03075* **2021**.
57. Zniyed, Y.; Nguyen, T.P.; et al. Enhanced network compression through tensor decompositions and pruning. *IEEE Transactions on Neural Networks and Learning Systems* **2024**.
58. Ding, B.; Qin, C.; Liu, L.; Chia, Y.K.; Li, B.; Joty, S.; Bing, L. Is GPT-3 a Good Data Annotator? In Proceedings of the ACL (1). Association for Computational Linguistics, 2023, pp. 11173–11195.
59. OpenAI. GPT-4V(ision) System Card. 2023.
60. Chiang, W.L.; Li, Z.; Lin, Z.; Sheng, Y.; Wu, Z.; Zhang, H.; Zheng, L.; Zhuang, S.; Zhuang, Y.; Gonzalez, J.E.; et al. Vicuna: An Open-Source Chatbot Impressing GPT-4 with 90%* ChatGPT Quality, 2023.
61. Zhao, B.; Wu, B.; He, M.; Huang, T. SVIT: Scaling up Visual Instruction Tuning, 2023, [\[arXiv:cs.CV/2307.04087\]](https://arxiv.org/abs/2307.04087).

62. Abdine, H.; Chatzianastasis, M.; Bouyioukos, C.; Vazirgiannis, M. Prot2Text: Multimodal Protein's Function Generation with GNNs and Transformers. In Proceedings of the Deep Generative Models for Health Workshop NeurIPS 2023, 2023.
63. Yang, Z.; Cherian, S.; Vucetic, S. Data Augmentation for Radiology Report Simplification. In Proceedings of the Findings of the Association for Computational Linguistics: EACL 2023; Vlachos, A.; Augenstein, I., Eds., Dubrovnik, Croatia, 2023; pp. 1922–1932. <https://doi.org/10.18653/v1/2023.findings-eacl.144>.
64. Xi, Y.; Liu, W.; Lin, J.; Cai, X.; Zhu, H.; Zhu, J.; Chen, B.; Tang, R.; Zhang, W.; Zhang, R.; et al. Towards Open-World Recommendation with Knowledge Augmentation from Large Language Models, 2023, [\[arXiv:cs.IR/2306.10933\]](https://arxiv.org/abs/2306.10933).
65. Luo, R.; Zhao, Z.; Yang, M.; Dong, J.; Li, D.; Lu, P.; Wang, T.; Hu, L.; Qiu, M.; Wei, Z. Valley: Video Assistant with Large Language model Enhanced ability, 2023, [\[arXiv:cs.CV/2306.07207\]](https://arxiv.org/abs/2306.07207).
66. Ouyang, L.; Wu, J.; Jiang, X.; Almeida, D.; Wainwright, C.; Mishkin, P.; Zhang, C.; Agarwal, S.; Slama, K.; Ray, A.; et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems* **2022**, *35*, 27730–27744.
67. Glaese, A.; McAleese, N.; Trebacz, M.; Aslanides, J.; Firoiu, V.; Ewalds, T.; Rauh, M.; Weidinger, L.; Chadwick, M.; Thacker, P.; et al. Improving alignment of dialogue agents via targeted human judgements. *arXiv preprint arXiv:2209.14375* **2022**.
68. Sun, Z. A short survey of viewing large language models in legal aspect. *arXiv preprint arXiv:2303.09136* **2023**.
69. Wei, W.; Ren, X.; Tang, J.; Wang, Q.; Su, L.; Cheng, S.; Wang, J.; Yin, D.; Huang, C. LLMRec: Large Language Models with Graph Augmentation for Recommendation, 2024, [\[arXiv:cs.IR/2311.00423\]](https://arxiv.org/abs/2311.00423).
70. Liu, Q.; Chen, N.; Sakai, T.; Wu, X.M. ONCE: Boosting Content-based Recommendation with Both Open- and Closed-source Large Language Models, 2023, [\[arXiv:cs.IR/2305.06566\]](https://arxiv.org/abs/2305.06566).
71. Kim, Y.; Rush, A.M. Sequence-level knowledge distillation. *arXiv preprint arXiv:1606.07947* **2016**.
72. Song, F.; Yu, B.; Li, M.; Yu, H.; Huang, F.; Li, Y.; Wang, H. Preference ranking optimization for human alignment. *arXiv preprint arXiv:2306.17492* **2023**.
73. Solaiman, I.; Dennison, C. Process for adapting language models to society (palms) with values-targeted datasets. *Advances in Neural Information Processing Systems* **2021**, *34*, 5861–5873.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.