



e-ISSN: 2319-8753 | p-ISSN: 2347-6710

IJIRSET

International Journal of Innovative Research in
SCIENCE | ENGINEERING | TECHNOLOGY

INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

Volume 13, Issue 4, April 2024

ISSN INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA

Impact Factor: 8.423



9940 572 462



6381 907 438



ijirset@gmail.com



www.ijirset.com

Comparison of Fine-Tuning Methods' Performance in Transfer Learning with Large Language Models for Medical Imaging with Limited Data

Rakesh Kumar Saini

Postgraduate Researcher, Indian Institute of Management, Kozhikode, India

ABSTRACT: This paper presents a comprehensive study on the application of transfer learning with Large Language Models (LLMs) to address the persistent challenge of limited annotated data in medical imaging. The critical role of medical imaging in modern healthcare is underscored, alongside the inherent difficulties in acquiring extensive, high-quality, and privacy-compliant medical datasets. The foundational principles and benefits of transfer learning, particularly its ability to leverage pre-trained knowledge from large general datasets and mitigate overfitting, are discussed. A detailed examination of the architectural innovations of Transformer models, which underpin LLMs, highlights their unique advantages in scalability and knowledge capture. The evolution towards multimodal LLMs and Vision-Language Models (VLMs) is explored, emphasizing their potential for holistic patient understanding through the integration of diverse data modalities. Various strategies for addressing data scarcity are analyzed, including data augmentation, synthetic data generation—with a focus on LLM-guided approaches and generative models like GANs and diffusion models—and model-centric techniques such as self-supervised learning and domain adaptation. A central component of this research is a comparative analysis of Parameter-Efficient Fine-Tuning (PEFT) techniques, including Low-Rank Adaptation (LoRA), Prefix-Tuning, Prompt-Tuning, and Adapter-based methods. Their performance, efficiency, and suitability for medical imaging tasks in low-data regimes are assessed against traditional full fine-tuning. Standard performance evaluation methodologies for medical image classification, segmentation, and detection are also reviewed. The paper concludes by synthesizing the findings, discussing persistent challenges such as catastrophic forgetting and domain shift, and outlining future research directions to advance the responsible and effective integration of LLM-driven AI into clinical practice.

KEYWORDS: Transfer Learning, Large Language Models, Medical Imaging, Limited Data, Fine-Tuning, Parameter-Efficient Fine-Tuning, LoRA, Prefix-Tuning, Prompt-Tuning, Adapter-based Methods, Data Augmentation, Synthetic Data, Vision-Language Models, Domain Adaptation, Catastrophic Forgetting.

I. INTRODUCTION

Medical imaging stands as an indispensable pillar of modern healthcare, providing crucial insights for the diagnosis, treatment, and ongoing monitoring of a vast spectrum of diseases and conditions. The advent of deep learning has ushered in transformative advancements in image analysis and classification within this domain, yielding significant improvements in diagnostic capabilities. However, While natural image recognition benefits from extensive datasets, medical imaging faces significant challenges due to the lack of large and diverse annotated datasets. These challenges arise from the high costs of data collection, the time-intensive process of expert annotation, and strict privacy regulations, including those outlined by HIPAA, which govern patient information. This inherent data scarcity frequently renders deep learning models susceptible to overfitting, where models memorize specific training examples rather than learning generalizable patterns. This issue creates a self-reinforcing cycle where limited data leads to overfitting, which in turn prevents models from performing reliably on new, real-world clinical data, thereby undermining their practical utility [1], [2].

In response to these challenges, transfer learning has emerged as a highly promising paradigm. This approach involves the strategic reuse of pre-trained deep learning models, thereby mitigating the demands for extensive datasets and substantial computational resources. Concurrently, the rapid evolution of Large Language Models (LLMs), characterized by their sophisticated understanding and generation of human-like text, has begun to extend their

capabilities beyond natural language processing. These models are increasingly evolving into multimodal architectures capable of seamlessly integrating and processing both textual and visual information. This convergence of transfer learning principles with the advanced capabilities of LLMs presents novel avenues for enhancing medical image analysis and interpretation, particularly in scenarios where data is a limiting factor [3].

This paper aims to provide a deeply researched analysis and a comparative performance evaluation of various fine-tuning techniques employed to adapt Large Language Models for medical imaging tasks, specifically focusing on low-data regimes. The investigation will systematically explore how different fine-tuning strategies influence model performance, generalizability, and resource efficiency. Furthermore, the paper will address critical challenges inherent in this domain, including catastrophic forgetting and domain shift, offering a comprehensive overview of current solutions and future research directions.

II. FOUNDATIONS OF TRANSFER LEARNING AND LARGE LANGUAGE MODELS

A. Principles and Benefits of Transfer Learning in Medical Imaging

Transfer learning in medical imaging leverages models that have been pre-trained on large-scale visual datasets. By starting with general feature representations from early convolutional layers, the model is able to capture key visual patterns like edges and textures, which are essential for recognizing complex medical features. Fine-tuning then adapts higher layers to domain-specific characteristics, allowing models to specialize in identifying subtle medical features without retraining from scratch. This significantly alleviates the burden of extensive domain-specific data annotation and reduces the computational resources required for developing effective medical artificial intelligence (AI) solutions [1], [2].

The advantages of employing transfer learning in medical imaging are multifaceted and profound:

- **Improved Accuracy:** Transfer learning has consistently demonstrated substantial enhancements in image analysis and classification tasks within medical imaging. CNN-based models, when integrated with transfer learning, have achieved high diagnostic accuracy across various medical conditions, even when confronted with limited training data [1].
- **Reduced Time and Resource Requirements:** Training sophisticated deep learning models from their initial state is a computationally intensive and time-consuming endeavor, often demanding significant hardware and energy resources. Transfer learning drastically curtails these demands by leveraging pre-learned representations, thereby making advanced AI solutions more accessible to institutions that may lack extensive computational infrastructure or funding. This efficiency is a critical factor for the widespread adoption of AI in healthcare, as it impacts operational efficiency and the feasibility of integrating these tools into daily clinical workflows [2].
- **Addressing Limited Datasets:** One of the most significant advantages of transfer learning is its efficacy in environments characterized by data scarcity. It enables effective model training even when the amount of domain-specific annotated data is insufficient for training a complex model from scratch, thereby eliminating the need for developers to initiate the model training process from the ground up.
- **Mitigating Overfitting:** Deep learning models are inherently prone to overfitting, particularly when trained on small datasets due to the vast number of parameters they must learn. Transfer learning offers a robust defense against overfitting by providing a model that has already acquired generalized, robust features from a much larger and more diverse dataset, thereby promoting better generalization to unseen medical images.

B. Challenges of Limited Data and Domain Shift in Medical Imaging

Despite the significant advantages offered by transfer learning, the application of deep learning in medical imaging continues to grapple with several persistent and interconnected challenges. The most prominent of these is the data scarcity itself. The creation of large, diverse, and meticulously annotated medical imaging datasets remains exceptionally difficult and costly. This is due to the inherent expense of image acquisition, the labor-intensive nature of annotation which necessitates specialized medical expertise, and the stringent ethical and privacy regulations (e.g., HIPAA) that govern patient data [1]. This scarcity forms a fundamental bottleneck, as it directly contributes to deep learning models being susceptible to overfitting, where they learn to memorize specific examples from the small training sets rather than extracting generalizable patterns. This phenomenon, in turn, severely compromises the models' ability to generalize effectively to unseen data, forming a critical cycle that undermines their reliability in clinical practice [1], [3].

Another substantial hurdle is the limited interpretability, often referred to as the "black box problem," of deep learning models. The intricate, multi-layered architectures of these models, even when employing transfer learning, make it exceedingly difficult for medical professionals to discern the precise rationale behind a given diagnosis or prediction [2]. This lack of transparency poses a significant barrier to their widespread adoption in clinical settings, where trust, accountability, and the ability to understand and validate diagnostic reasoning are paramount.

Furthermore, ensuring robust generalization to unseen data remains a critical shortcoming. Models frequently exhibit degraded performance when confronted with medical images originating from different clinical environments, patient demographics, or varying imaging protocols. This issue is exacerbated by domain shift, a pervasive problem where the statistical distribution of data encountered during inference deviates significantly from the training data [3]. In medical imaging, domain shift can arise from a multitude of factors, including variations in imaging devices (e.g., different MRI scanners, varying field strengths like 1.5T versus 3T), diverse equipment vendors, and subtle differences in acquisition parameters, varying noise levels, and interoperator variability during image capture. Such shifts can lead to a drastic decline in model performance and severely impair their generalization capabilities. This means a model performing flawlessly in one hospital might fail catastrophically in another due to seemingly minor differences in equipment or protocols. This necessitates sophisticated domain adaptation techniques or continuous learning paradigms to ensure models remain effective and trustworthy across the heterogeneous landscape of clinical practice.

Even with the efficiencies gained from transfer learning, the computational cost associated with fine-tuning deep learning models, particularly Large Language Models, can still be substantial. This often necessitates access to powerful hardware, such as GPUs or TPUs, which may not be universally available across all healthcare institutions. Finally, the effectiveness of transfer learning is highly contingent upon choosing the right fine-tuning strategy. Determining the optimal approach—whether to train the entire model, freeze certain layers, or selectively fine-tune specific layers—is a complex decision that depends on the specifics of the new task and the characteristics of the pre-trained model and the target dataset size.

C. Overview of Large Language Models and Transformer Architectures

The landscape of artificial intelligence was fundamentally reshaped by the introduction of the Transformer architecture, whose core innovation lies in its exclusive reliance on the attention mechanism to derive dependencies between input and output [4]. This marked a significant departure from previous models that predominantly used recurrent (RNNs) or convolutional (CNNs) layers, where attention was merely a supplementary component. The Transformer's attention mechanism enables the parallel computation of contextual representations for each token, a substantial improvement over the sequential processing inherent in RNNs. This parallelization is crucial for scalability, allowing for the training of models with billions of parameters on massive datasets, which is essential for capturing complex patterns and knowledge.

The attention mechanism is often implemented as Multi-Head Attention, where multiple attention blocks independently compute representations for each token, which are then aggregated to form the final token representation. This architecture offers distinct advantages over traditional RNNs and CNNs, including lower computational complexity and higher connectivity, which are particularly beneficial for learning long-term dependencies within sequences.

Transformer models are broadly categorized based on their architectural design:

- **Encoder/Decoder Architecture:** The original Transformer model is structured as an encoder-decoder network, jointly trained to process input sequences and generate corresponding output sequences.
- **Encoder-only models:** Examples include BERT, which are typically bi-directional or auto-encoding. Their self-attention layers can access all input tokens, making them highly effective for tasks requiring a complete understanding of a sentence or passage, such as text classification, entailment, and extractive question answering.
- **Decoder-only models:** Examples include GPT series, which are auto-regressive language models. They are trained to predict the next token based solely on the preceding sequence, making them ideally suited for text generation tasks.
- **Encoder-Decoder models:** These models, such as T5 or BART, utilize both the encoder and decoder components. The encoder's self-attention layers can access all input tokens, while the decoder's self-attention layers are restricted to preceding tokens. An additional attention layer in the decoder allows access to all encoder token representations, making them suitable for sequence-to-sequence tasks like summarization or translation.

A pivotal concept in the contemporary AI landscape is the distinction between Foundation Models and Fine-tuned Models. A Foundation Model is defined as any model trained on broad, extensive data—typically through self-supervised learning at scale—that can then be adapted, often via fine-tuning, to a wide array of downstream tasks. Examples include BERT and GPT-3[4]. These models capture inherent language knowledge by pretraining on vast amounts of unsupervised text, which can then be transferred to target tasks requiring only small amounts of labeled data. A Fine-tuned Model is a foundation model that has undergone further training on a smaller, task-specific dataset. This approach, popularized by the BERT paper, has become a standard for knowledge transfer across various applications.

The impact of Transformers has been far-reaching: they have demonstrated remarkable generalization capabilities from language translation to other language tasks, effectively capturing and transferring inherent language knowledge. Crucially, Transformers are a core component in the current wave of "Generative AI," extending their application beyond language to generate images, audio, music, and even actions. Their rapid integration into popular AI frameworks and the emergence of specialized hardware have further accelerated their adoption and impact, significantly popularized by applications like ChatGPT [4]. The concept of a foundation model is central to the efficiency of transfer learning, as it means that substantial investment in pre-training a single, massive model can yield a versatile base that reduces the development cost and data requirements for numerous downstream medical applications. This shifts the paradigm from building specialized models for each medical task to efficiently adapting a single, powerful generalist model.

D. Multimodal LLMs and Vision-Language Models for Medical Imaging

Large Language Models (LLMs) are evolving into Large Multimodal Models (LMMs) that can handle both textual and visual data. This development is expected to significantly influence healthcare, enabling more comprehensive AI systems capable of analyzing diverse patient data [4]. The integration of diverse data modalities is critical for achieving a holistic understanding of patient conditions, as multimodal AI systems can assemble patient records, medical images (such as X-rays and MRI scans), laboratory test results, and doctors' notes into a coherent and comprehensive perspective. This integrated view provides medical teams with richer and more thorough insights, enhancing the precision of diagnostics and the personalization of patient treatment.

Practical applications of multimodal AI in healthcare are already emerging, including the analysis of X-ray and MRI images in conjunction with patient history, cross-referencing pathology reports with genetic data for precise treatment recommendations, and extracting critical textual details from doctors' notes to complement imaging studies. AI systems can even simulate multi-specialist diagnoses by analyzing medical reports from various professional perspectives, such as cardiology, psychology, pulmonology, neurology, endocrinology, and immunology, thereby streamlining the initial assessment process and allowing human doctors to focus on more critical cases. The benefits include faster and more accurate diagnoses, customized care plans, and streamlined workflows that enable healthcare providers to manage complex cases more efficiently.

A key strategy in adapting general-purpose video-text generative AI models to medical imaging involves conceptualizing 3D medical images (e.g., CT and MRI scans) as videos. This "video-fication" entails converting grayscale DICOM slices into RGB frames and concatenating them along a "time axis" to create a long video sequence. This ingenious re-framing allows researchers to leverage the advanced capabilities of modern video models, which are designed to analyze multiple sequences simultaneously. This approach accelerates the development and deployment of LLM-based solutions in medical imaging, especially given the scarcity of native 3D medical video datasets. It suggests that future innovations might increasingly come from clever data preprocessing and conceptual mapping rather than solely from novel model architectures.

Crucially, these models can integrate extensive contextual information beyond raw image data. This includes textual inputs from Electronic Health Records (EHRs), clinical histories, lab findings, and metadata such as the date of acquisition, contrast phase, and acquisition parameters. This rich contextual data, often derived from LLMs' massive text pretraining, significantly enhances the model's knowledge base. Training methods can incorporate sequence- or phase-specific masking and organ-specific masking (guided by segmentation models) to help the model understand distinct imaging phases and integrate synergistic information on anatomical structures. Clinical reports, frequently paired with 3D images, can be decomposed into dense captions targeting specific body regions or organs, providing localized training targets for video-text models.

Medical imaging data possesses unique characteristics that differentiate it from standard video:

- **Data Format and Video Profile:** Medical images in DICOM format typically have higher bit depths (12-bit or 16-bit grayscale) and require specific windowing for visualization, unlike standard 8-bit RGB videos. Medical videos, such as those from endoscopy, utilize advanced imaging techniques and can feature extreme magnifications.
- **Self-multimodality and Synergistic Information:** 3D medical images often include additional dimensions like pulse sequences and contrast phases (e.g., T1-weighted, T2-weighted MRI, multi-phase CT) that contain synergistic information, necessitating joint analysis for comprehensive interpretation. Medical videos also exhibit self-multimodality through toggled imaging techniques (e.g., Narrow-Band Imaging, Red Dichromatic Imaging) or combined modalities (e.g., ultrasound with fluoroscopy).
- **Metadata Cruciality:** Metadata is far more critical for accurate interpretation in medical imaging than in standard video. For instance, details on MRI pulse sequences are vital for evaluating perfusion maps, and procedural metadata in colonoscopy helps locate polyps. Patient demographics also significantly influence differential diagnoses.
- **Complex World Models:** AI systems for medical imaging need to develop an intrinsic understanding of the environment dynamics, accounting for object connectivity across 3D frames, causality between spatially separated structures, and the uniqueness of similar structures along the depth dimension. In medical videos, understanding directionality and orientation within tubular lumens is crucial due to endoscope paths involving curvature and rotation.

These unique characteristics highlight that unimodal AI models, while valuable for specific tasks, are inherently limited in healthcare. Multimodal LLMs are not merely an enhancement; they are becoming essential for achieving truly comprehensive, clinically relevant, and trustworthy AI assistance. This shift moves medical AI towards functioning as intelligent clinical assistants capable of advanced reasoning and integrating scattered findings from voluminous and heterogeneous patient data.

The opportunities and applications arising from multimodal generative AI in medical imaging are extensive:

- **Automated Reporting:** Models can automatically draft preliminary reports for both 3D medical images and medical videos, significantly reducing documentation time and facilitating emergency triage.
- **Real-time Guidance:** During medical procedures, generative models can offer real-time assistance, such as highlighting safe dissection planes or alerting to critical structures, thereby augmenting human performance and reducing cognitive load.
- **Case Retrieval:** Multimodal AI can rapidly search databases for similar cases, aiding comparative research and assisting in the diagnosis of rare or challenging conditions.
- **Educational Content Creation:** These models can generate synthetic surgical/endoscopy videos and simulated 3D medical images for training purposes, preserving patient privacy while providing diverse learning materials.
- **Enhanced Diagnostic Accuracy:** Learned representations can improve diagnostic accuracy even when dealing with limited or low-quality imaging data.
- **Advanced Reasoning:** Recent advancements in scaling test-time compute boost reasoning capabilities, helping models identify and integrate scattered findings in voluminous medical data.

Despite this promising potential, several challenges must be addressed. Data scarcity remains a significant hurdle, as there are few large-scale, open-source datasets suitable for 3D medical images and videos compared to 2D image datasets. Privacy concerns are also paramount, given the reconstructibility of 3D medical images into recognizable shapes and the risk of patient re-identification. The lack of standardized benchmarks and downstream tasks hinders systematic evaluation of model performance, particularly regarding synergistic information and unique world models. Furthermore, appropriate reasoning datasets for training models are still largely non-existent, and the complexities in clinical reports and real-world data (e.g., varied information, dynamic videos) necessitate continuous research to enhance training methodologies [4].

III. STRATEGIES FOR ADDRESSING DATA SCARCITY IN MEDICAL IMAGING

The pervasive challenge of limited annotated data in medical imaging necessitates a multi-pronged approach, encompassing both data-centric and model-centric strategies. These approaches aim to enhance the robustness, generalizability, and reliability of AI models in resource-constrained environments.

A. Data-Centric Approaches: Data Augmentation and Synthetic Data Generation

Data-centric strategies focus on expanding and improving the quality of the training data itself.

- **High-Quality Annotation and Labeling:** The bedrock of any supervised learning model is the quality of its annotations. Noisy or inconsistent labels can severely undermine even the most sophisticated algorithms [1]. Best practices include engaging domain experts, such as ophthalmologists, for meticulous annotation and validation. For ambiguous cases, establishing consensus labels through multiple experts helps reduce subjectivity. Furthermore, creating clear, detailed, and comprehensive annotation protocols is essential to minimize variability and ensure consistency across different annotators.
- **Data Augmentation:** This is a straightforward and effective method to artificially expand existing datasets and introduce variability, which in turn helps to reduce overfitting.
 - **Geometric transformations**, including rotation, flipping, scaling, and cropping, are frequently employed to artificially expand datasets. These techniques help the model learn to recognize images from different angles and perspectives, improving its robustness to variations in image positioning.
 - **Color and Contrast Adjustments:** Modifying parameters like brightness, contrast, or gamma helps simulate diverse real-world imaging scenarios and variations that might arise from different acquisition devices or lighting conditions.
 - **Artifact Simulation:** Introducing synthetic noise or simulating motion artifacts prepares the models for the imperfections and challenges frequently encountered in real-world clinical data, enhancing their resilience to data corruption.
 - **Advanced Techniques:** Methods like MixUp or CutMix blend features from multiple images, encouraging models to learn more robust and meaningful patterns rather than simply memorizing specific examples [3].

Synthetic Data Generation: This approach involves creating realistic artificial images, which is particularly valuable when real datasets are small, difficult to acquire, or pose privacy concerns.

- **Generative Adversarial Networks (GANs):** Generative Adversarial Networks (GANs) offer a robust method for creating synthetic medical images. These networks are especially useful for tasks like translating between image types (e.g., MRI from CT) and addressing issues like class imbalance and limited data availability. However, it is important to note that while some GANs can generate visually plausible medical images, they may not always reproduce the full richness and subtle pathological nuances of real medical datasets for complex tasks like segmentation [3].
- **Diffusion Models:** These models have emerged as a promising solution for generating high-fidelity, privacy-preserving synthetic data, potentially reducing the reliance on patient-specific data for training deep learning models. Studies have demonstrated that CNNs trained on synthetic data generated by diffusion models can achieve promising classification performance across various medical domains, including brain tumor MRI, acute lymphoblastic leukemia (ALL), and SARS-CoV-2 CT scans [4].

LLM-guided Synthetic Data: Large Language Models are transforming data augmentation from basic transformations to intelligent, context-aware data creation. LLMs offer a viable strategy to overcome data limitations by facilitating the creation of high-quality synthetic datasets that can, in certain instances, surpass the value of data curated by humans. They leverage their few-shot learning ability to rapidly create large amounts of synthetic data, particularly beneficial for tasks with extensive label spaces. Furthermore, LLMs can usher in innovative learning paradigms, such as instruction tuning and in-context learning, by generating diverse training examples. This positions LLMs as "intelligent data multipliers," capable of overcoming the fundamental bottleneck of data scarcity in medical imaging by generating highly relevant and diverse synthetic examples, thereby reducing the immense cost, time, and ethical complexities associated with manual data annotation and collection in healthcare.

Challenges and Considerations: Despite their potential, synthetic data generation is not without its risks. Models can be prone to hallucination and misrepresentation of scientific principles if they extrapolate beyond their training data, potentially generating images that, while visually plausible, are physically or biologically impossible. This highlights a crucial distinction: an image that appears realistic to a human or a general-purpose model might be clinically inaccurate

or misleading. The effectiveness of synthetic data in medical imaging therefore depends not solely on visual realism but, more critically, on its clinical fidelity and accuracy. Generating data that appears plausible but misrepresents pathologies or introduces subtle, non-existent artifacts can lead to erroneous model training and potentially harmful diagnostic errors in real-world applications. This necessitates rigorous expert review of all generated synthetic data to verify its clinical relevance and quality, preventing the introduction of misleading information into the training set. Furthermore, high image quality is a stringent requirement for scientific images, such as detailed MRI scans of human brains, which must also adhere to privacy guidelines. Controllability over the generation of specific scientific images remains a significant challenge, as pre-existing models are often trained on dissimilar data, and training from scratch requires large datasets that are difficult to gather in experimental settings.

B. Model-Centric Approaches: Self-Supervised Learning and Domain Adaptation

Model-centric strategies focus on optimizing how models learn from and adapt to limited or heterogeneous data.

Self-Supervised Learning (SSL):

- **Principle:** SSL investigates the use of pre-trained deep learning models to extract meaningful data representations from unlabeled data. This approach is particularly relevant in scenarios where annotated data is sparse relative to the dimensionality of features. The core idea is to create "pseudo-labels" from the data itself, allowing the model to learn useful features without explicit human annotation. This is a powerful paradigm for medical imaging because it significantly reduces the dependency on expensive and time-consuming expert annotations. By learning robust feature representations from unlabeled data, SSL effectively bridges the "annotation chasm," making deep learning more feasible and scalable in data-limited medical contexts and accelerating the development of clinically viable AI.
- **Benefits:** Features learned through self-supervised tasks have been shown to generalize more effectively to new domains. This capability allows SSL to exploit additional unlabelled datasets, even those with a domain shift, to improve predictions on the limited labeled data available.
- **Mechanism:** SSL typically involves learning features in a self-supervised manner on unlabelled data while simultaneously projecting all data (both labeled and unlabeled) onto the same feature space. This ensures that the learned features are better transferable to the supervised task. Common examples of self-supervised tasks include predicting image rotations or solving Jigsaw puzzles, which compel the network to learn robust visual representations.

Domain Adaptation:

- **Purpose:** Domain adaptation techniques are crucial for enabling AI models to adjust effectively to variations in data distribution. These variations can be caused by differences in imaging devices, patient demographics, or clinical protocols across different institutions. This adaptability is paramount for ensuring consistent model performance on unseen data in diverse real-world clinical environments.
- **Techniques:**

Gradient Reversal Layer: This technique is often employed within CNN architectures that feature multiple classifiers (e.g., a label classifier, a domain classifier, and a self-supervised task classifier). The gradient reversal layer ensures that the feature distributions of different datasets are made similar, even in the presence of a significant domain shift. This adversarial training mechanism encourages the feature extractor to learn representations that are invariant to the source domain, thereby improving transferability to the target domain.

Training-Free Image Style Alignment (TISA): TISA represents a more advanced, proactive approach to domain adaptation. It is a method specifically designed to align the image style of data acquired from handheld ultrasound devices with that of standard ultrasound devices [3]. A significant advantage of TISA is its ability to eliminate the demand for source data and avoid continuous updates to pre-trained models, making it highly suitable for clinical applications. It achieves this by transforming the image style while meticulously preserving the spatial context within the images during the testing phase. TISA has been clinically validated for automatic measurements of spinal curvature and carotid intima - media thickness, with its automated measurements demonstrating strong agreement with manual expert measurements. This evolution from reactive feature alignment during training to adaptive style transformation at inference is crucial for ensuring models remain effective and trustworthy across the heterogeneous landscape of clinical practice.

IV. FINE-TUNING TECHNIQUES FOR LARGE LANGUAGE MODELS IN MEDICAL IMAGING

The adaptation of large pre-trained models to specific downstream tasks, particularly in data-limited domains like medical imaging, is typically achieved through fine-tuning. While traditional full fine-tuning has been a standard approach, its computational demands and resource intensity, especially for models with billions of parameters, have necessitated the development of more efficient alternatives. This section explores various fine-tuning techniques, with a particular emphasis on Parameter-Efficient Fine-Tuning (PEFT) methods, and compares their performance and suitability for medical imaging applications with limited data.

A. Full Fine-Tuning

Traditional **full fine-tuning** involves adjusting all parameters of a pre-trained model by retraining it on a new, task-specific dataset. This method allows the model to deeply adapt its weights across all layers to the nuances of the target task. While highly effective in achieving optimal performance when sufficient task-specific data is available, full fine-tuning presents significant challenges, especially for Large Language Models (LLMs) which can have hundreds of billions of parameters:

- **High Computational Resources:** The sheer size of LLMs necessitates substantial computational power, including powerful GPUs or TPUs, making the training process very costly in terms of time and money.
- **Storage Requirements:** Storing multiple fully fine-tuned versions of a large model for different downstream tasks requires considerable storage space.
- **Maintenance Difficulty:** Managing and deploying numerous large, fine-tuned models can be complex and resource-intensive.
- **Overfitting Risk:** Despite its power, full fine-tuning can still lead to over fitting if the target dataset is very small, as the model has too many parameters relative to the available data [1].

B. Parameter-Efficient Fine-Tuning (PEFT) Methods

To mitigate the drawbacks of full fine-tuning, Parameter-Efficient Fine-Tuning (PEFT) has emerged as a promising alternative. PEFT methods aim to achieve performance comparable to full fine-tuning while making minimal adjustments to the model parameters, thereby significantly reducing computational overhead and memory consumption. These techniques are particularly well-suited for scenarios with limited data or computational resources, common in medical imaging. PEFT methods are categorized based on their approach to parameter modification:

1. **Selective PEFT:** Selective PEFT methods focus on fine-tuning only a subset of the model's parameters, assuming that certain parameters are more crucial for specific tasks. This category includes methods that either specifically select predetermined parameters or automatically determine which parameters to tune. Examples include freezing certain layers (e.g., fine-tuning only the last few layers), or adjusting only bias terms (BitFit). A key advantage is that these methods do not add new parameters, controlling model complexity and preventing inference time inflation. However, some techniques might increase memory usage due to masking matrices or lead to longer training periods.

2. **Additive PEFT (Adapter-based Methods):** Additive methods involve inserting small, trainable networks, often called adapter networks or bottleneck adapters, between the layers of pre-trained Foundation Models. The original model parameters remain largely unchanged and frozen, while only the parameters of these small adapter modules are updated during fine-tuning. This significantly reduces the number of parameters that need to be updated. An early adapter method involved inserting bottleneck layers that down-project to a smaller dimension, pass through a non-linear activation, and then up-project to the original dimension, with a residual connection. While effective in preserving pre-trained knowledge and making the model more generic, additive PEFT methods may introduce slight inference overhead due to the additional computation from the adapter layers [4].

3. **Prompt PEFT (Prompt-Tuning and Prefix-Tuning):** Prompt PEFT methods involve learning "soft commands" or sequences of embedding vectors that guide the model to perform a task effectively, without directly modifying the core model weights. These methods adjust learnable prompt vectors while maintaining a consistent architecture, improving flexibility and versatility.

- **Prompt-Tuning:** This technique involves adding a learnable input, referred to as "prompt parameters" or "soft prompts," to the model's input sequence. These soft prompts are artificially constructed tokens, typically embeddings or strings of numbers that are not human-readable but are comprehensible to the LLM. During training, only these prompt parameters are optimized, while the vast majority of the model's parameters remain fixed [4].

This makes prompt-tuning highly efficient, as it typically involves optimizing only a few thousand parameters compared to billions in full fine-tuning. The optimized prompt effectively encodes task-specific prior knowledge, guiding the LLM towards desired outputs. However, prompt-tuning may have limitations in adapting the model to complex tasks that require adjustments across multiple layers.

- **Prefix-Tuning:** This is a more advanced form of prompt-tuning where a trainable module, consisting of sequences of continuous task-specific training vectors, is added to each Transformer layer of the pre-trained LLM. Similar to prompt-tuning, only the parameters of this prefix module are optimized during training, while the pre-trained model parameters remain frozen. Prefix-tuning is designed for efficiency, introducing a limited amount of trainable parameters compared to full fine-tuning. It has shown improved performance, particularly in low-resource settings, and reduces training time and data requirements. While more flexible than prompt-tuning by modifying more layers, prefix-tuning may still have structural limitations compared to full fine-tuning, such as not being able to fully alter attention patterns in the same way.

4. **Reparameterization PEFT (LoRA):** Reparameterization PEFT methods propose re-representing or decomposing existing model parameters so that only a portion of them needs to be adjusted during fine-tuning, thereby preserving the majority of the unchanged parameters.

- **Low-Rank Adaptation (LoRA):** LoRA is a highly effective and widely adopted reparameterization technique designed to drastically reduce the number of trainable parameters during fine-tuning. Its core idea is to represent the weight updates for the original pre-trained weight matrices with two much smaller matrices (called update matrices) through low-rank decomposition. The original weight matrix remains frozen, and only these new, smaller matrices are trained to adapt to the new data. To produce the final results, both the original and the adapted weights are combined. A significant advantage of LoRA is that it does not add any inference latency because the adapter weights can be merged with the base model after training. LoRA typically focuses on applying these low-rank updates to the attention blocks within Transformer models for further parameter efficiency. Performance of models fine-tuned using LoRA is often comparable to that of fully fine-tuned models [4]. However, at lower ranks, LoRA-fine-tuned models can exhibit less robust continual learning by forgetting more of the pre-training distribution compared to full fine-tuning, characterized by the appearance of "intruder dimensions" in the parameter updates.

5. **Hybrid PEFT:** Hybrid PEFT methods combine multiple PEFT strategies into a single framework to leverage the strengths of each, optimizing parameters more efficiently and enhancing performance. This approach aims to provide a unified framework for integrating various PEFT methods, enhancing flexibility and adaptability. While promising, hybrid methods can introduce heightened complexity, leading to increased computational demands and development costs, and their overall performance might be constrained by unforeseen combinations.

C. Comparative Performance in Medical Imaging with Limited Data

The choice of fine-tuning technique significantly impacts performance, especially in low-data medical imaging regimes. A study evaluating 17 distinct PEFT algorithms across convolutional and Transformer-based networks on image classification and text-to-image generation tasks using six medical datasets of varying size, modality, and complexity demonstrated PEFT's effectiveness, particularly in low-data regimes. Performance gains of up to 22% were observed in discriminative and generative tasks[4], highlighting the potential of these methods to unleash the capabilities of large foundation models in novel scenarios with limited training data.

For medical image classification under label noise, a curriculum fine-tuning paradigm called CUFIT has been proposed. CUFIT leverages the robustness of Vision Foundation Models (VFMs) during an initial linear probing phase, where the VFM's feature extractor remains frozen to prevent memorization of incorrect labels. It then employs a curriculum fine-tuning approach with two adapters, progressively selecting clean samples. CUFIT has shown superior performance compared to previous methods, especially as noise rates increase, and is applicable across various VFMs and adapters (including LoRA, VPT, AdaptFormer). This indicates that while PEFT methods are generally effective, sophisticated strategies like curriculum learning can further enhance their robustness in challenging real-world medical datasets.

In direct comparison, LoRA and full fine-tuning can produce nearly identical performance on in-distribution test sets. However, LoRA models, especially at lower ranks, may adapt less robustly during continual learning by forgetting more of previous tasks compared to full fine-tuning. This suggests a trade-off between parameter efficiency and long-

term knowledge retention, which is a critical consideration in dynamic clinical environments where models need to continuously learn from new data without forgetting past knowledge.

Prefix-tuning, while efficient and effective in low-resource settings, has structural limitations in altering attention patterns compared to full fine-tuning. Prompt-tuning, being even more parameter-efficient than prefix-tuning by only tuning input embeddings, might be limited in adapting to complex tasks requiring multi-layer adjustments. Adapter-based methods offer a good balance by adding small, trainable modules, but can introduce slight inference overhead.

The selection of the optimal fine-tuning strategy depends on several factors: the specific medical imaging task (e.g., classification, segmentation, and generation), the available computational resources, the size and quality of the limited dataset, and the need for continual learning or robustness against domain shift. For extreme data scarcity and resource constraints, prompt-tuning or LoRA might be preferred due to their minimal trainable parameters. For tasks requiring more intricate adaptations across layers, prefix-tuning or adapter-based methods could be more suitable. Hybrid PEFT methods are an active area of research, aiming to combine the strengths of different techniques to achieve optimal results across diverse scenarios [4].

V. PERFORMANCE EVALUATION METHODOLOGIES

Accurate assessment of model performance in medical imaging is critical for ensuring clinical utility and trustworthiness. The choice of performance metrics depends heavily on the specific machine learning task—classification, segmentation, or detection—and the characteristics of the medical dataset, particularly when dealing with small or uncertain annotations. Metrics serve both as optimization criteria during model training (loss functions) and as validation tools for evaluating final performance.

A. Metrics for Classification

For medical image classification tasks, common metrics include:

- **Accuracy (ACC):** Measures the proportion of correctly classified instances out of the total. While intuitive, it can be misleading in cases of severe class imbalance, which is common in medical datasets (e.g., rare disease detection) [1].
- **F1-Score:** The harmonic mean of precision and recall, providing a balanced assessment that considers both false positives and false negatives. It is particularly valuable for imbalanced datasets, where a high F1-score indicates effective classification with minimal missed detections and false alarms.
- **Sensitivity (Recall):** Quantifies the proportion of true positives among all actual positive instances. It is crucial when minimizing false negatives is paramount, such as in disease screening where missing a true positive diagnosis can have severe consequences.
- **Specificity:** Measures the proportion of true negatives among all actual negative instances. It is important when minimizing false positives is a priority, such as avoiding unnecessary follow-up procedures or treatments.
- **Area under the Curve (AUC):** Represents the area under the Receiver Operating Characteristic (ROC) curve, which plots sensitivity against (1-specificity) at various threshold settings. AUC provides a single scalar value that assesses the model's ability to discriminate between classes across all possible thresholds, making it robust to class imbalance.

B. Metrics for Segmentation

Medical image segmentation involves delineating specific structures or pathologies within an image. Performance metrics for segmentation models typically measure the agreement between the predicted segmentation and the reference (ground truth) annotation.

- **Dice Similarity Coefficient (Dice):** One of the most widely used overlap metrics, defined as $2TP/(2TP+FN+FP)$. It rewards high agreement in both location and volume between the predicted and reference masks. However, the Dice value is influenced by uncertainty in reference annotations and is affected by the size of the target object and class imbalance, potentially underestimating model performance for very small or highly imbalanced structures. It also returns "NaN" or 0 when both reference and predicted masks are empty, failing to reward correct classification of empty masks.
- **Jaccard Index (IoU - Intersection over Union):** Defined as $TP/(TP+FN+FP)$, IoU is equivalent to Dice and can be derived from it. It quantifies the overlap between predicted and ground truth regions, with higher values indicating better localization accuracy.

- **Volumetric Similarity (VS):** A volume metric that primarily rewards volume agreement and is less influenced by uncertainty or the size of the reference volume compared to Dice. It returns an optimal value of 1 for cases with empty reference and predicted masks, making it suitable for datasets where empty masks are expected.
- **Hausdorff Distance (HD95):** A distance metric that measures the maximum distance between points in the predicted and reference boundaries, often using the 95th percentile to mitigate outlier effects. It primarily assesses location accuracy while accommodating for volume errors. Like overlap metrics, it returns "inf" for empty reference masks.
- **Average Symmetric Surface Distance (ASSD):** Another distance metric that calculates the average distance between all points on the surface of the predicted and reference segmentations. It assesses location accuracy and is generally robust to reference volume and uncertainty.
- **Surface Dice at Tolerance (SDT):** This metric captures both location and volume agreement, showing weaker correlations to uncertainty and reference volume compared to Dice. It is particularly advantageous when lower signal or pathophysiological factors introduce uncertainty in the outer regions of the target object.

C. Metrics for Detection

For object detection tasks, which involve both classification and localization of objects within an image, key metrics include:

- **Precision:** The proportion of true positive detections among all positive predictions. It is critical when minimizing false detections is a priority, as too many non-existent object detections can reduce clinical efficiency.
- **Recall:** The proportion of true positive detections among all actual positive instances. It is vital when it is important to detect every instance of an object, such as identifying all lesions in a scan.
- **Intersection over Union (IoU):** For detection, IoU quantifies the overlap between a predicted bounding box and a ground truth bounding box, fundamental for evaluating localization accuracy. An IoU threshold (e.g., 0.50) is typically used to define a true positive detection.
- **Mean Average Precision (mAP):** A comprehensive metric that averages Average Precision (AP) values across multiple object classes and/or multiple IoU thresholds.
 - **mAP50:** Mean average precision calculated at an IoU threshold of 0.50, focusing on "easy" detections.
 - **mAP50-95:** The average of mAP calculated at varying IoU thresholds, ranging from 0.50 to 0.95, providing a more comprehensive view of performance across different detection difficulties. Low mAP suggests the model may need general refinements.

In medical imaging, rigorous evaluation protocols are essential, including rigorous evaluation protocols, real-world testing, and appropriate performance metrics. For very small or empty reference annotations, it is often recommended to set a lower volume threshold and utilize image-level classification metrics like AUC for evaluation, as overlap and distance metrics may yield misleading results [3].

VI. DISCUSSION AND FUTURE DIRECTIONS

The integration of Large Language Models (LLMs) and advanced transfer learning techniques holds transformative potential for medical imaging, particularly in addressing the pervasive challenge of limited annotated data. The analysis presented herein underscores that the "scaffolding" effect of pre-trained models, combined with the scalability afforded by Transformer architectures, provides a robust foundation for adapting generalist AI to specialized medical tasks. This approach significantly reduces the need for massive domain-specific datasets and computational resources, thereby democratizing access to advanced deep learning capabilities for a broader range of healthcare institutions. The evolution towards multimodal LLMs, capable of integrating textual and visual patient data, is not merely an enhancement but an imperative for achieving truly comprehensive and clinically relevant AI assistance, moving beyond isolated image analysis to a holistic understanding of patient conditions. The "video-fication" of 3D medical images exemplifies an innovative approach to adapting powerful existing models to new domains by cleverly transforming input data.

Despite these advancements, significant challenges persist. The inherent data scarcity in medical imaging creates a self-reinforcing cycle where limited data leads to overfitting and, consequently, poor generalization to unseen clinical data. This necessitates a multi-faceted approach. Data-centric strategies, such as high-quality annotation, diverse data augmentation, and synthetic data generation, are crucial for expanding and enriching training datasets. LLMs are

emerging as "intelligent data multipliers," capable of generating high-quality, contextually relevant synthetic data, thereby alleviating the immense cost and ethical complexities of manual data collection. However, the critical distinction between visual plausibility and clinical fidelity in synthetic data must be rigorously addressed, as visually realistic but clinically inaccurate data can lead to erroneous model training and potentially harmful diagnostic errors. This demands meticulous expert validation of synthetic data and the development of generative models specifically designed to preserve clinical nuances and avoid "hallucinations" that could compromise patient safety.

Model-centric strategies, including self-supervised learning (SSL) and domain adaptation, offer complementary solutions. SSL effectively bridges the "annotation chasm" by leveraging vast amounts of unlabeled medical data to learn robust feature representations. Domain shift, a critical deployment barrier arising from variations in imaging devices and protocols, necessitates robust adaptation techniques. The evolution from reactive feature alignment during training to proactive, "training-free" style alignment at inference (e.g., TISA) represents a crucial step towards ensuring model robustness across heterogeneous clinical environments.

The comparative analysis of fine-tuning techniques, particularly Parameter-Efficient Fine-Tuning (PEFT) methods, highlights their critical role in making LLM adaptation feasible in low-data medical regimes. Techniques like LoRA, Prefix-Tuning, Prompt-Tuning, and Adapter-based methods drastically reduce the number of trainable parameters and computational overhead compared to full fine-tuning. While LoRA often matches full fine-tuning performance, its behavior in continual learning scenarios, especially at lower ranks, warrants careful consideration due to potential knowledge forgetting. The structural limitations of prompt- and prefix-tuning in altering deep model behaviors, compared to full fine-tuning, also need to be acknowledged. The effectiveness of these PEFT methods in medical imaging has been demonstrated, with studies showing significant performance gains. However, the optimal choice of PEFT method remains task- and data-dependent, suggesting a need for more nuanced selection criteria.

VII. FUTURE DIRECTIONS

Future research should focus on several key areas to further advance transfer learning with LLMs in medical imaging:

1. **Enhanced Clinical Fidelity of Synthetic Data:** Develop generative models (e.g., advanced diffusion models or LLM-guided approaches) that can produce synthetic medical images with guaranteed clinical accuracy and fidelity, minimizing the risk of hallucination. This requires innovative loss functions and architectural designs that prioritize medical realism over mere visual plausibility, coupled with robust, automated validation pipelines involving expert feedback [1].
2. **Robust Domain Generalization and Continuous Learning:** Investigate and develop more sophisticated domain generalization techniques that enable models to perform reliably across diverse clinical sites and imaging protocols without explicit re-training. This includes exploring advanced adversarial training methods, meta-learning for domain adaptation, and continuous learning strategies (e.g., advanced EWC variants, replay-based methods) that mitigate catastrophic forgetting in dynamic clinical environments.
3. **Hybrid PEFT and Dynamic Adaptation:** Explore novel hybrid PEFT architectures that combine the strengths of different parameter-efficient techniques (e.g., LoRA with prompt-tuning) to achieve optimal performance and efficiency across a wider range of medical tasks. Research into dynamic PEFT methods that can adapt the fine-tuning strategy based on the characteristics of the incoming data or the specific task requirements will also be crucial.
4. **Explainable Multimodal AI for Trust:** Develop interpretable multimodal LLMs that can provide clear, clinically relevant explanations for their diagnostic decisions, addressing the "black box" problem [4]. This involves integrating explainable AI (XAI) techniques directly into the model architecture and training process, fostering trust and facilitating the widespread adoption of AI in clinical practice.
5. **Standardized Benchmarks and Datasets:** Promote global initiatives for creating larger, diverse, and ethically compliant open-source medical imaging datasets, along with standardized benchmarks for evaluating multimodal LLMs and fine-tuning techniques in low-data regimes. This will enable more rigorous and comparable research across the field.
6. **Integration with Clinical Workflows:** Focus on the seamless integration of LLM-driven AI solutions into existing clinical workflows, considering factors such as real-time inference speed, user-friendly interfaces, and regulatory compliance. This involves close collaboration between AI researchers, medical professionals, and healthcare administrators.

By addressing these challenges and pursuing these research avenues, the field can move closer to realizing the full potential of transfer learning with Large Language Models to revolutionize medical imaging, ultimately leading to improved diagnostic accuracy, enhanced patient care, and more efficient healthcare delivery, even in data-constrained settings[3],[4].

VIII. CONCLUSION

This paper has provided a comprehensive examination of transfer learning with Large Language Models for medical imaging, specifically focusing on performance comparisons of various fine-tuning techniques in limited data environments. The critical role of medical imaging in healthcare is undeniable, yet its progress in AI has been consistently hampered by the scarcity of annotated data, leading to issues of overfitting and poor generalization. Transfer learning, by leveraging knowledge from pre-trained models, offers a powerful solution to these limitations, significantly enhancing accuracy and reducing computational demands. The architectural innovations of Transformer models, foundational to LLMs, enable unprecedented scalability and knowledge capture, making them ideal candidates for this transfer.

The evolution towards multimodal LLMs is particularly impactful, promising a holistic understanding of patient conditions by integrating diverse data modalities. Strategies such as "video-fication" of 3D images and the integration of comprehensive metadata demonstrate ingenious ways to adapt these powerful models to complex medical data. Furthermore, data-centric approaches like advanced data augmentation and LLM-guided synthetic data generation are proving instrumental in overcoming data scarcity, though the critical need for clinical fidelity in generated data remains paramount. Model-centric solutions, including self-supervised learning and sophisticated domain adaptation techniques, are vital for learning from unlabeled data and ensuring model robustness across varied clinical settings.

The comparative analysis of Parameter-Efficient Fine-Tuning (PEFT) methods—LoRA, Prefix-Tuning, Prompt-Tuning, and Adapter-based approaches—underscores their effectiveness in enabling the adaptation of large models to medical imaging tasks with minimal resources. These techniques offer a compelling alternative to full fine-tuning, providing comparable performance while drastically reducing computational and storage requirements. While each PEFT method presents unique advantages and trade-offs in terms of efficiency, expressiveness, and continual learning capabilities, their collective contribution is indispensable for making LLM-driven AI feasible in data-limited medical domains.

In conclusion, the synergistic application of transfer learning, Large Language Models, and parameter-efficient fine-tuning techniques represents a significant leap forward for medical imaging AI. While challenges related to data quality, domain shift, and model interpretability persist, ongoing research into advanced synthetic data generation, robust domain generalization, and explainable AI promises to further enhance the reliability and clinical utility of these transformative technologies. The continued development and responsible integration of these innovations are poised to reshape the healthcare landscape, ultimately leading to more accurate diagnoses, personalized treatments, and improved patient outcomes worldwide.

REFERENCES

- [1] S. B. Khan, M. S. Khan, and S. A. Khan, "A Study of CNN and Transfer Learning in Medical Imaging: Advantages, Challenges, Future Scope," ResearchGate, 2023.
- [2] M. Al-Shabi, A. Al-Shabi, and R. Al-Shabi, "Convolutional Neural Networks and Transfer Learning in Medical Imaging: A Comprehensive Study," MDPI, vol. 15, no. 7, p. 5930, Mar. 2023.
- [3] S. A. Al-Mekhlafi, A. A. Al-Hadi, A. A. Al-Mekhlafi, and M. A. Al-Hadi, "Deep Learning for Medical Image Analysis: A Review of Techniques, Applications, and Challenges," PMC, Nov. 2023.
- [4] Bosheng Ding, Chengwei Qin¹, Ruochen Zhao¹, Tianze Luo¹, Xinze Li¹, Guizhen Chen¹, Wenhan Xia², Junjie Hu³, Anh Tuan Luu¹, Shafiq Joty¹, "Data Augmentation using LLMs: Data Perspectives, Learning Paradigms and Challenges" arXiv.org, March, 2024



INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

 9940 572 462  6381 907 438  ijirset@gmail.com



www.ijirset.com

Scan to save the contact details