

Prompt Mechanisms in Medical Imaging: A Comprehensive Survey

Hao Yang^{1,+}, Xinlong Liang^{2,3,+}, Zhang Li^{4,+}, Yue Sun¹, Zheyu Hu^{5,6}, Xinghe Xie¹, Behdad Dashtbozorg^{2,7}, Jincheng Huang¹, Shiwei Zhu¹, Luyi Han^{2,3}, Jiong Zhang⁸, Shanshan Wang⁹, Ritse Mann^{2,3,*}, Qifeng Yu^{4,*}, and Tao Tan^{1,*}

¹Faculty of Applied Sciences, Macao Polytechnic University, Rua de Luis Gonzaga Gomes, Macao, China

²Netherlands Cancer Institute, Department of Radiology, Amsterdam, 1066 CX, The Netherlands

³Radboud University Medical Centre, Department of Radiology and Nuclear Medicine, Nijmegen, 6525 GA, The Netherlands

⁴College of Aerospace Science and Engineering, National University of Defense Technology, Changsha, China

⁵Medical Department of Breast Cancer, Hunan Cancer Hospital, Changsha, China

⁶Medical Department of Breast Cancer, the Affiliated Cancer Hospital of Xiangya School of Medicine, Central South University, Changsha, China.

⁷Faculty of Biomedical Engineering, Eindhoven University of Technology, Antonie van Leeuwenhoek, Plesmanlaan 121, Eindhoven, Noord Brabant, NL

⁸Laboratory of Advanced Theranostic Materials and Technology, University of Chinese Academy of Sciences, China

⁹Paul C. Lauterbur Research Center for Biomedical Imaging, Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, China

*corresponding authors: Ritse.Mann@radboudumc.nl; yuqifeng@nudt.edu.cn; taotanjs@gmail.com

+these authors contributed equally to this work

ABSTRACT

Deep learning offers transformative potential in medical imaging, yet its clinical adoption is frequently hampered by challenges such as data scarcity, distribution shifts, and the need for robust task generalization. Prompt-based methodologies have emerged as a pivotal strategy to guide deep learning models, providing flexible, domain-specific adaptations that significantly enhance model performance and adaptability without extensive retraining. This systematic review critically examines the burgeoning landscape of prompt engineering in medical imaging. We dissect diverse prompt modalities, including textual instructions, visual prompts, and learnable embeddings, and analyze their integration for core tasks such as image generation, segmentation, and classification. Our synthesis reveals how these mechanisms improve task-specific outcomes by enhancing accuracy, robustness, and data efficiency and reducing reliance on manual feature engineering while fostering greater model interpretability by making the model's guidance explicit. Despite substantial advancements, we identify persistent challenges, particularly in prompt design optimization, data heterogeneity, and ensuring scalability for clinical deployment. Finally, this review outlines promising future trajectories, including advanced multimodal prompting and robust clinical integration, underscoring the critical role of prompt-driven AI in accelerating the revolution of diagnostics and personalized treatment planning in medicine.

Introduction

The integration of Artificial Intelligence (AI), intense learning models such as convolutional neural networks (CNNs), has revolutionized medical image analysis, producing remarkable advancements in image classification, segmentation, and generation¹⁻⁴. Despite these successes, clinical translation and widespread adoption of these models are often impeded by persistent challenges such as data scarcity and distribution shifts. Among these, distribution shifts between training datasets and various real-world clinical data present a significant obstacle, often leading to a degradation in model performance and reliability⁵. These shifts arise from inherent dataset biases, inter-scanner variability, and diverse image acquisition protocols. Consequently, improving model adaptability and robustness to such heterogeneous data distributions has become a significant research imperative in medical image processing⁶⁻¹³.

The recent advent of Foundation Models (FMs), large-scale models pre-trained on extensive and diverse datasets, underpinned by architectures such as the Transformer and encompassing architectures such as Large Language Models (LLMs), and crucially for this field, Vision-Language Models (VLMs), offers a transformative new paradigm for AI in medical imaging¹⁴. These FMs exhibit superior versatility and transfer learning capabilities, demonstrating immense potential in few-shot learning,

cross-modal information fusion, and robust performance in complex data-scarce scenarios. However, effectively harnessing the power of these generalist FMs for specialized medical tasks requires sophisticated adaptation strategies beyond traditional fine-tuning. In this context, prompt-based mechanisms have rapidly emerged as a key and powerful approach to steering FMs. Concurrently, the underlying principles of prompting-guiding model behavior with external contextual information have also been increasingly explored and adapted for a broader range of deep learning architectures in medical imaging, extending beyond large-scale FMs. These mechanisms are proving particularly effective in addressing the aforementioned distribution shifts, managing complex specialized medical knowledge, and unlocking the full potential of FM and other customized models for applications in precision medicine and personalized diagnostics.

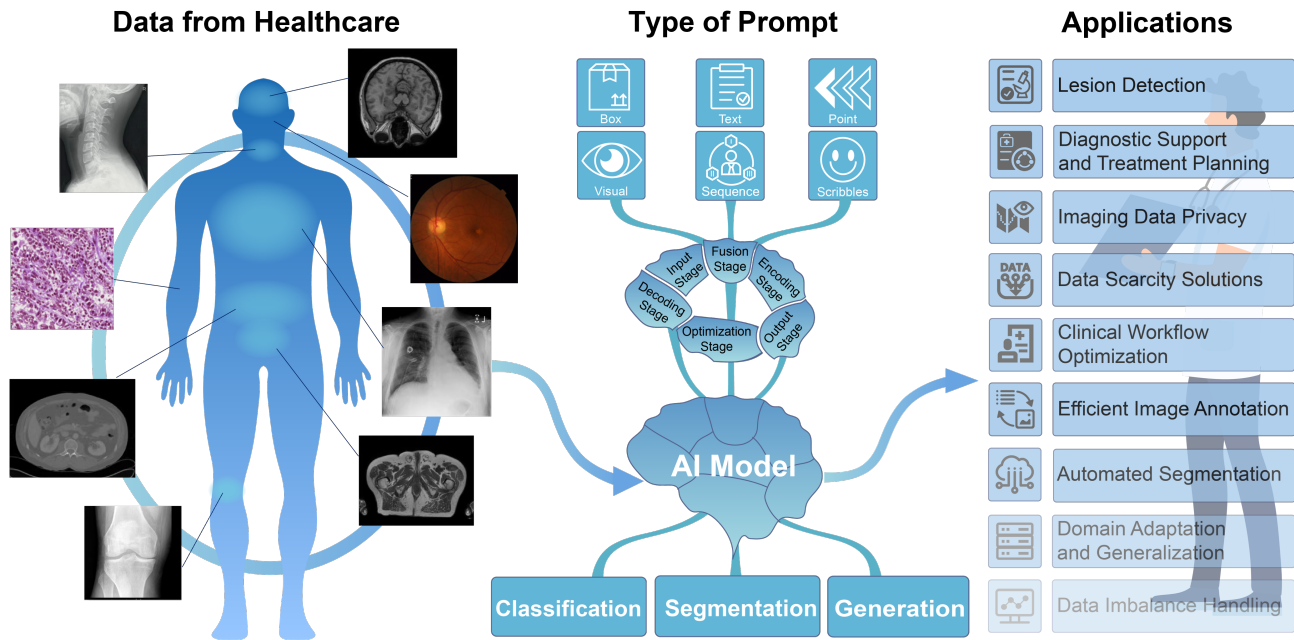


Figure 1. An overview of the framework for prompt-based AI in medical imaging. The process begins with various medical imaging data sources, which are processed by an AI model capable of accepting multiple prompt types (visual, text, sequence, etc.). By performing classification, segmentation, or generation tasks, the model ultimately serves a range of clinical applications to improve diagnostic efficiency and accuracy.

The core tenet of prompting is to guide and optimize a model's performance on specific downstream tasks by introducing targeted external contextual information. As illustrated in the general framework of figure 1, this process begins with various sources of medical imaging data. An AI model, designed to be receptive to multiple prompt types (visual, text, etc.), is then guided by these targeted inputs. The prompt leads the model to perform a specific task, such as classification, segmentation, or generation, ultimately serving a range of clinical applications by enhancing diagnostic accuracy and efficiency. This information can be encoded as textual instructions⁶, visual prompts (e.g., points, boxes)¹⁵, or learnable embeddings. This strategy aims to significantly augment the model's capacity to interpret and analyze medical images with high fidelity, especially when annotated data is scarce or expensive. A key technical advantage of prompting lies in its remarkable parameter efficiency: it often necessitates fine-tuning only a small subset of the model's parameters-or none at all in zero-shot prompting scenarios-to adapt the pre-trained model to new tasks, thereby obviating the need for substantial architectural modifications or complete retraining from scratch⁴. Furthermore, prompting offers exceptional flexibility and adaptability, with diverse guidance strategies enabling models to tackle a broad spectrum of medical imaging tasks and their associated challenges¹⁶. By facilitating the integration of additional knowledge and context, prompting can significantly enhance the generalization and robustness of the model, particularly when faced with limited sample sizes⁵⁻⁸. Recent advances in this growing field showcase increasingly sophisticated techniques, including multimodal prompts that fuse information from various channels (e.g., image and text) to tackle complex problems¹⁷, adaptive prompting algorithms that dynamically tailor prompt content based on specific task requirements and data characteristics¹⁶, and personalized prompting schemes that customize inputs for individual images to mitigate issues arising from sample distribution discrepancies, thereby enhancing overall system performance¹⁸.

Despite these promising developments and the rapid proliferation of research, the practical application and optimal design

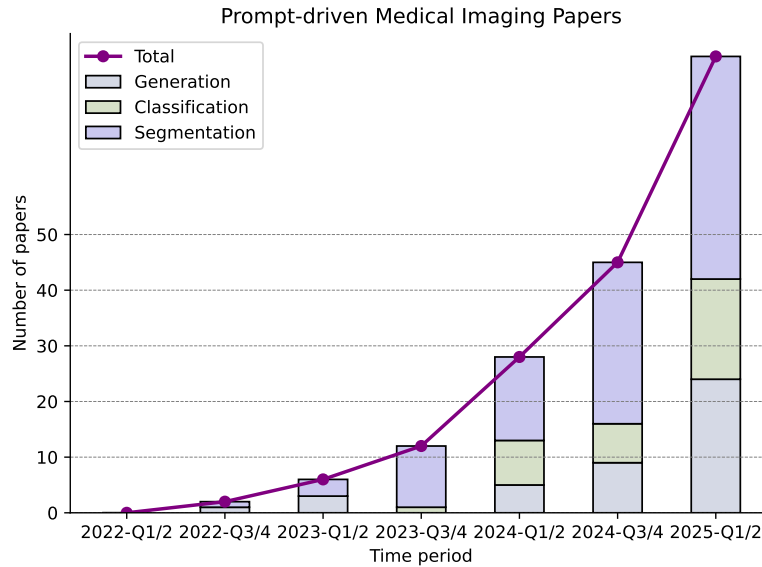


Figure 2. Rapid increase of the number of Prompt-driven Medical imaging papers. Generation, Classification, and Segmentation are the three main taxonomy categories introduced in this survey.

of prompt mechanisms in medical imaging are still confronted by several critical limitations deeply intertwined with the nature of prompting. While prompting is often positioned as a solution for data scarcity¹¹, the design and validation of high-quality prompts introduce their data-dependent challenges. Crafting effective textual prompts that capture nuanced clinical knowledge or generating representative learnable embeddings requires a sufficient volume of diverse and meticulously annotated data⁵⁻⁸, the availability of which is curtailed by stringent privacy regulations. The complexities of medical annotation further undermine the quality and reliability of prompts; the process demands specialized expertise and is prone to inter-observer variability. This subjectivity translates directly into inconsistent or ambiguous prompts-whether textual, visual, or learned-which can mislead the model and degrade performance, complicating the development of standardized and reproducible prompting protocols.

Crucially, the challenge of distribution shift⁵ extends beyond the model's core to the prompts themselves. A prompt that is effective for one data distribution (e.g., images from a specific scanner or institution) may fail dramatically on another, a phenomenon that can be described as prompt brittleness. This lack of robustness is a significant barrier, as prompts must be generalizable across different patient populations and the notable domain gaps among imaging modalities (e.g., CT, MRI, X-ray, Ultrasound). These multifaceted factors-spanning prompt design, quality, and robustness-collectively constrain the generalizability and clinical applicability of prompt-guided models, including FMs. This underscores the urgent need for research into more adaptive, robust, and systematically designed prompting strategies that can overcome these inherent obstacles.

Given the rapid advancements and the critical need for consolidated knowledge in this burgeoning area, this paper reviews prompt engineering tailored explicitly for medical imaging. The urgency and relevance of this review are underscored by the exponential growth in related academic publications, as depicted in figure 2. The data clearly shows a steep increase in the number of papers on prompt-driven medical imaging, with significant contributions across the primary tasks of generation, classification, and segmentation. This trend highlights the intense research interest in the field and necessitates a timely, structured synthesis of existing work. We aim to comprehensively analyze the diverse methodologies and widespread applications of prompt mechanisms across various model architectures. Our approach involves an extensive literature search to gather and synthesize relevant studies, followed by developing and applying a novel classification system. This system categorizes existing research into two primary dimensions: the core technologies underpinning prompt mechanisms (encompassing design, generation, integration strategies, and their synergy with transfer learning and multimodal learning) and their clinical application paradigms (evaluating effectiveness in classification, segmentation, generation, and their impact on diagnostic accuracy and treatment planning). This review summarizes the distinct advantages and inherent limitations of various prompting methods through a detailed, multi-faceted analysis of the collated literature. It critically evaluates their performance across multiple crucial dimensions, including accuracy, robustness, and data utilization efficiency. Our significant contributions are:

- **Systematic and Comprehensive Review:** Providing a structured and in-depth survey of the application and methodology of prompt mechanisms in medical imaging, integrating theoretical advancements with practical applications, and

underscoring their transformative potential in image generation, segmentation, and classification tasks.

- **Novel Classification Framework and Trend Analysis:** Introducing an innovative classification system that categorizes research into core technologies and clinical applications. This framework delineates current trends and innovations in prompt design, generation, and embedding techniques, offering clear and actionable directions for future research endeavors.
- **Identification of Current Challenges and Future Trajectories:** Delivering a thorough summary of the extant research challenges, including issues related to prompt quality, diversity, model scalability, and seamless integration into clinical workflows. Alongside this, we propose promising future development paths, such as advanced multimodal data integration, the pursuit of personalized medical solutions through tailored prompting, and breakthroughs in few-shot or zero-shot learning capabilities.

Theoretical Foundations and Taxonomies of Prompt Mechanisms

Prompt-based mechanisms represent a significant paradigm shift in guiding deep learning models, particularly for nuanced tasks in medical imaging where precision and domain-specific knowledge are paramount. Understanding these mechanisms' foundational principles and systematic classification is crucial for their practical design and application. This section outlines the core theoretical underpinnings of prompt engineering and presents a taxonomy for categorizing the diverse array of prompting strategies.

Theoretical Foundations

Core Theoretical Principles

Fundamentally, prompts act as conditioning signals that modulate neural network behavior, enabling precise output control, often without requiring extensive architectural modifications or complete retraining. From an information-theoretic point of view, prompts can be seen as injections of prior knowledge, such as domain-specific terminology, spatial cues, or task objectives, that constrain the hypothesis space, thus enhancing task-specific performance while aiming to maintain generalizability. The theoretical basis for prompts draws from established concepts such as transfer learning, where prompts facilitate knowledge transfer from pre-trained models; a Bayesian inference perspective, viewing prompts as informative priors guiding model learning; and manifold learning, where prompts help navigate the learned latent spaces of models towards desired outputs.

Key Operational Mechanisms

The efficacy of prompts comes from several key operational mechanisms. Attention modulation is primary, where prompts guide the model's internal attention to focus on task-relevant features while suppressing irrelevant information. Feature space transformation occurs when prompts induce changes in the model's representational geometry, projecting inputs into regions corresponding to desired outputs, often via parameter-efficient fine-tuning methods like prompt tuning, prefix tuning, or Low-Rank Adaptation, which apply low-rank updates to the model's weight matrices. Finally, conditional computation pathways may emerge in advanced architectures, where prompts activate specific computational subgraphs, enabling task-specialized processing.

Prompt Mechanisms in Medical Image Generation

Medical image generation, aiming to synthesize clinically and anatomically realistic images via deep learning, frequently necessitates meticulous conditioning to align outputs with specific clinical requirements or anatomical verisimilitude. Prompt mechanisms have become indispensable for achieving this alignment, effectively guiding generative models toward producing clinically relevant and structurally sound medical imagery. This subsection details the primary categories of prompts employed in this domain, highlighting how the general principles discussed previously are instantiated.

Text-Guided Synthesis: Leveraging Semantic Descriptions

Textual prompts are critical conditional controls in medical image generation, particularly for steering diffusion models and other generative architectures. These prompts, often derived from unstructured clinical narratives, professional medical reports (e.g., pathology reports¹⁹), or specific textual annotations detailing disease backgrounds⁶ or specific pathological features²⁰, encapsulate vital pathological details and contextual information. This semantic guidance ensures that generated images possess practical clinical significance and anatomical plausibility, which are crucial in fields like dermatology, chest X-ray synthesis, and 3D/4D volumetric data generation. Textual conditions furnish essential prior knowledge and empower developers to precisely direct the generation process by articulating specific disease characteristics. Text-driven synthesis typically relies on cross-attention mechanisms integrated within diffusion models to achieve high-fidelity medical image outputs.

Core Technical Pipeline for Textual Prompting From a technical standpoint, the core of this process often lies in powerful VLMs. These systems, exemplified by architectures such as BiomedCLIP²¹ and MedCLIP²² (as illustrated in figure 3), are fundamentally designed to align representations from medical images and textual descriptions, typically through a contrastive loss function that maximizes the similarity between corresponding image-text pairs. A key component of these architectures is a dedicated text encoder, which is responsible for parsing specialized medical terminology (e.g., "prurigo nodularis" or "squamous cell carcinoma"). To achieve high fidelity, these text encoders are often based on domain-specific language models such as PubMedBERT²³, ClinicalBERT²⁴, or other specialized Medical BERT variants²⁰. The primary function of these encoders, whether as part of a larger VLM or used independently, is to transform the text into high-dimensional embedding vectors that capture its rich semantic content. Subsequently, these embeddings are integrated into the generative model, typically through its layers of cross-attention^{25–28}. Within diffusion models, for instance, these textual embeddings play a pivotal role during the iterative denoising process^{2,6,29–35}, guiding the network to construct images that faithfully align with the provided textual descriptions. Researchers can precisely manipulate the generated image characteristics through detailed textual prompts by enabling fine-grained control over aspects such as lesion size, color, or anatomical location.

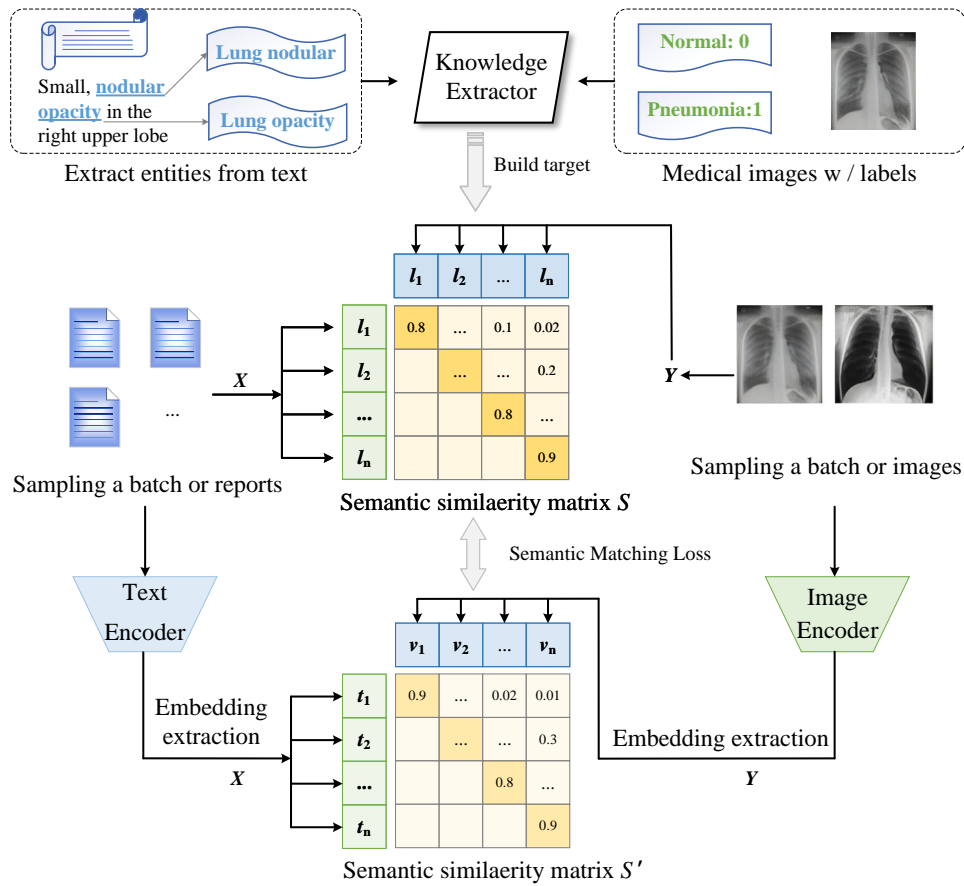


Figure 3. The MedCLIP²² framework for vision-language pre-training. It aligns medical image and text representations by training encoders to match a fine-grained semantic similarity matrix (bottom) derived from medical knowledge extraction (top).

Strategies for Optimizing Text-Conditioned Generation Various advanced strategies and auxiliary modules are often incorporated to enhance textual conditional control further and ensure the domain-specificity and realism of generated images. For example, hierarchical text prompt systems can be constructed using GPT-based summarizers¹⁹ or localized lesion description modules, facilitating the development of more robust disease representations within the generative model and improving the fidelity of images synthesized for rare or subtle lesions. Concurrently, many approaches integrate sophisticated alignment and consistency enforcement techniques. Some generation pipelines employ a pre-alignment step³⁵, where text tokens are first matched with corresponding disease-related regions in latent feature maps before proceeding with diffusion or GAN synthesis.

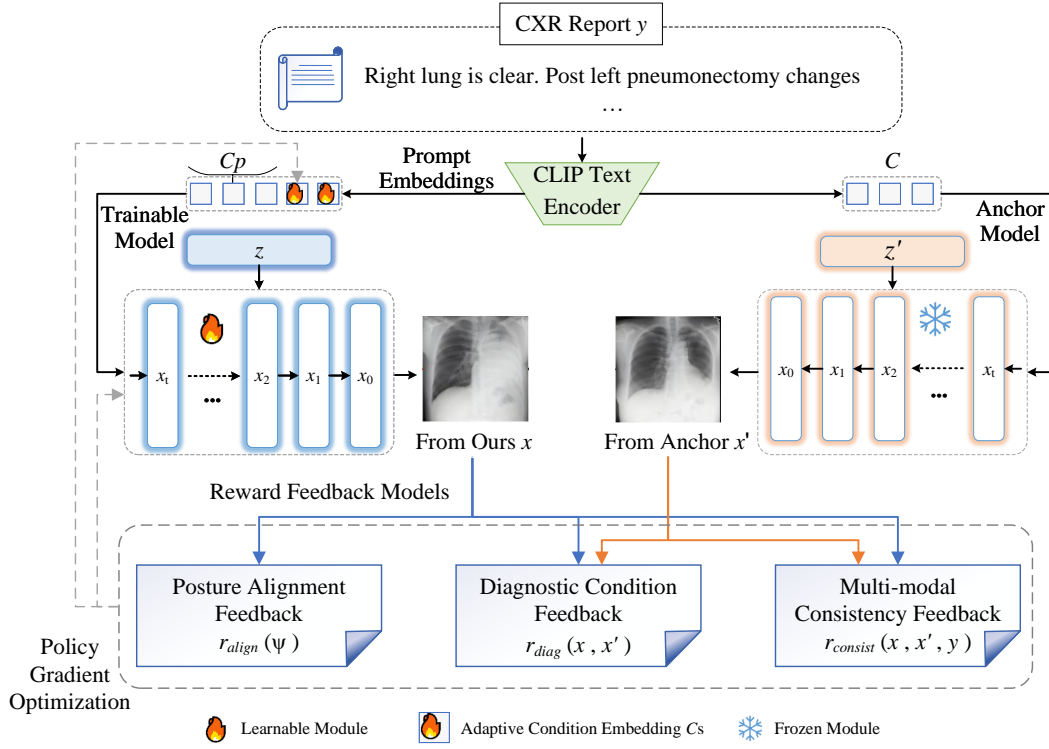


Figure 4. The CXRL framework³⁴, which utilizes textual reports as prompts to condition CXR generation. A key feature is using a reward feedback loop to refine the learnable prompt embeddings, ensuring alignment between the generated image and the diagnostic information in the report.

Other methodologies leverage domain adaptation techniques³¹ or concepts such as Textual Inversion². Textual Inversion, for instance, allows the model to learn a new pseudo-word embedding that represents a novel medical concept from just a few image examples, which can then be used in textual prompts to generate precise corresponding pathological features. The synergy between appropriate text encoders, advanced attention mechanisms, and these optimization strategies significantly elevates text-guided medical image synthesis’s realism, diversity, and controllability.

Non-Textual and Structurally-Aware Prompts in Generation

In addition to semantic guidance from text, the generation of medical images often benefits from prompts that provide structural, sequential, or other forms of explicit prior information. These non-textual or structurally-aware prompts enable finer-grained and more multidimensional control over the synthesis process, guiding the model from anatomical, temporal, or biological perspectives. Such prompts not only enhance the realism of image details and the consistency of anatomical structures but also help address clinical challenges like data scarcity and inter-sequence or inter-patient variability.

Anatomical Prompts for Structural Fidelity Anatomical prompts directly leverage key structural information, such as the delineations of lung lobes, airways, or blood vessels. This information is typically extracted from real CT or MRI images using pre-trained segmentation tools (e.g., lungmask³⁶, NaviAirway³⁷, and TotalSegmentator³⁸). During the generation process, these anatomical maps or constraints are injected into the generator, often through multichannel concatenation with latent representations²⁰ at various generator stages, such as along the input noise or at intermediate layers of a U-Net-like architecture. The generator is then tasked with producing the target intensity image (e.g., a single-channel CT slice) and ensuring consistency with the provided anatomical information. This imposes explicit structural constraints during both low-resolution generation and high-resolution super-resolution stages, effectively suppressing structural hallucinations and ensuring a higher degree of anatomical consistency and realistic detail in the synthesized images.

Sequence Prompts for Multi-Contrast MRI Synthesis Particularly relevant for MRI, sequence prompts guide the synthesis of specific image contrasts (e.g., T1-weighted, T1Gd, T2-weighted, FLAIR). These are typically represented using one-hot encodings or learnable embedding vectors that indicate the target MRI sequence. These prompts are injected as dynamic

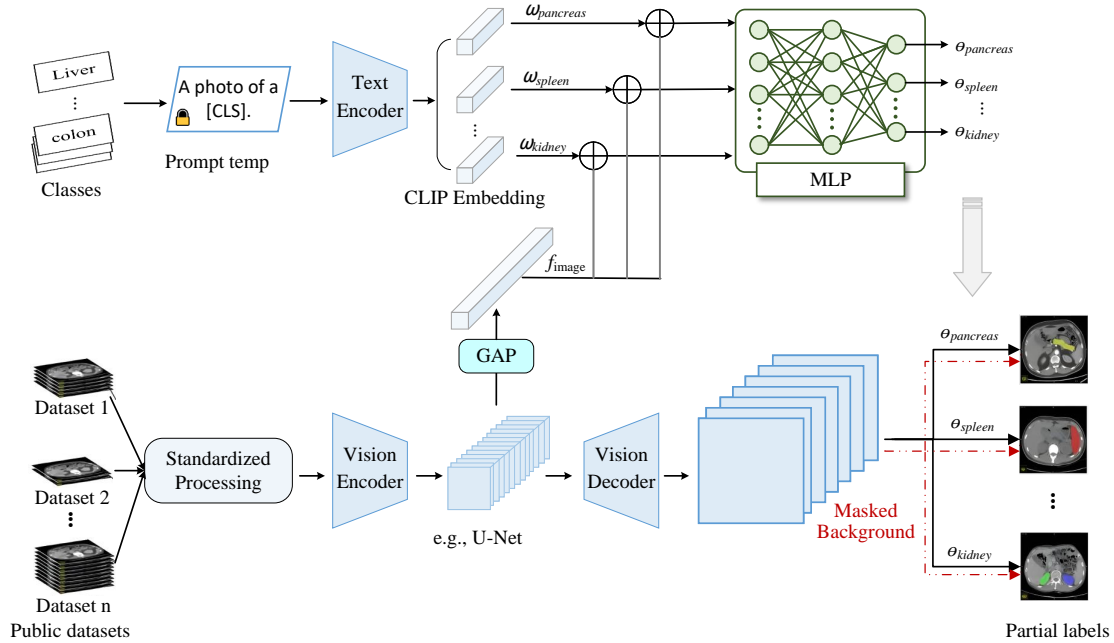


Figure 5. A CLIP-driven universal segmentation model⁷ where text prompts, generated from class names, are used to dynamically steer the network towards segmenting specific organs and tumors.

conditions into the generative network. For example, in³⁹, sequence prompts, in conjunction with structural features extracted by an encoder, drive a dynamic decoder. This allows a shared discrete latent representation to be mapped to the specific target sequence, facilitating cross-sequence generation from a common underlying anatomical representation. In another approach⁴⁰, predefined sequence codes are transformed into control vectors via a multilayer perceptron (MLP). These vectors then dynamically modulate a shared weight library within the HyperConv layers of a HyperDecoder module. This mechanism uses a smaller network (the hyper-network) to generate the weights (convolutional kernels) for the leading network (the HyperDecoder) based on the sequence prompt, allowing for dynamic, sequence-specific convolutions. This precisely generates the convolutional kernels specific to each target sequence while capturing shared (redundant) and unique (distinct) information across different MR sequences. Such mechanisms improve generative models' robustness and semantic expressiveness, especially under unsupervised or data-scarce conditions, and provide adequate support for handling missing clinical data or synthesizing entire multi-sequence MRI protocols.

Multi-Variable Prompts for Conditioned Phenotype Generation Multi-variable prompts have shown significant utility in generating images conditioned on specific patient characteristics or quantitative biomarkers, such as brain MRI generation. As shown in³, multiple normalized variables, including patient age, sex, ventricle volume, or normalized brain volume-can be injected into the latent space of a generative model, often via cascading them with latent vectors or using cross-attention mechanisms to allow the model to condition its output on these continuous or categorical variables at different levels of the generative process. This allows for dynamic control over the generation process, allowing the synthesis of brain images that reflect specific demographic or pathological states. Such methods exhibit strong controllability for directed generation (e.g., simulating aging effects). They can show excellent extrapolation capabilities, producing plausible image structures even for combinations of variables not explicitly seen during training.

Label Prompts for Genotype-Phenotype Association in Synthesis Label prompts can convert discrete categorical information, such as genotypic data (e.g., IDH mutation status, 1p19q co-deletion status in gliomas), into conditional vectors. These vectors are then injected into the conditioning mechanism of generative models, often by modulating the time or class embeddings within diffusion models⁴¹, to guide the synthesis process. This allows the network to learn and capture subtle morphological differences in images associated explicitly with these genotypes, influencing the denoising and reconstruction stages. This approach has proven effective for pathological image synthesis, producing high-quality images aligned with clinical or molecular characteristics. It offers considerable value for targeted data augmentation, especially for rare genetic

subtypes, and for visually exploring genotype-phenotype correlations.

The sophisticated integration of these diverse prompt mechanisms at various stages of generative models is crucial. It enhances the structural consistency and realistic detail of synthesized medical images and significantly strengthens the fine-grained controllability of the generation process, catering to a wide array of complex clinical conditions and research inquiries.

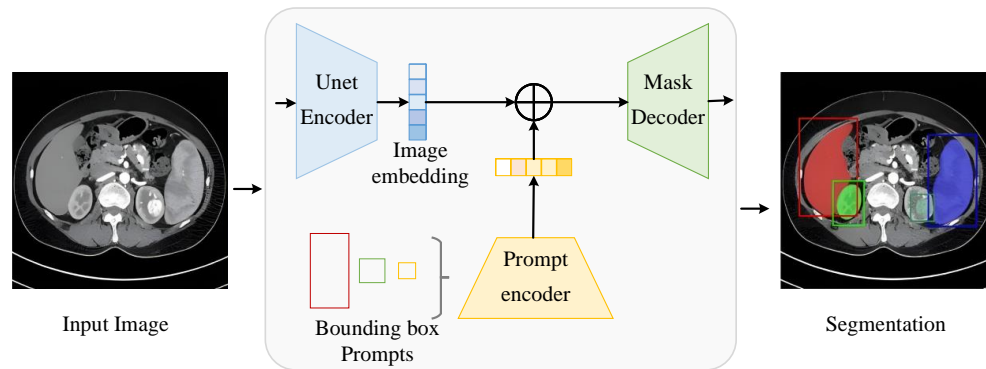


Figure 6. The architecture of MedSAM⁴², which utilizes bounding boxes as prompts to guide the segmentation process. The prompts are encoded and fused with the image embedding, directing the mask decoder to generate precise segmentations for the specified target regions.

Prompt Mechanisms in Medical Image Segmentation

In medical image segmentation, prompts play an essential role by guiding models to accurately focus on specific anatomical regions or pathological structures of interest, ultimately enhancing segmentation accuracy, robustness, and efficiency. Integrated at various stages of the model architecture, from initial input conditioning to feature extraction, attention modulation, and output refinement, prompts significantly influence segmentation outcomes. The effectiveness of a given prompt type is contingent upon the specific segmentation task, the characteristics of the available data (e.g., scarcity, annotation quality), and the nature of user interaction required or desired. This section classifies and examines the principal categories of prompts employed in contemporary medical image segmentation research. It focuses on their integration mode, impact on model performance, and strategic placement within the segmentation pipeline.

Visual Prompts for Interactive and Targeted Segmentation

Visual prompts have emerged as a potent tool for improving the accuracy and robustness of medical image segmentation, particularly valuable for interactive or semi-automated segmentation workflows. By providing explicit spatial prior information about the target regions, visual prompts direct the model's attention and facilitate precise delineation, especially in scenarios involving complex anatomical structures, low-contrast boundaries, or ambiguous image features. These prompts are pivotal across multiple stages of the segmentation model:

Input Stage: Initial Spatial Guidance Users typically provide simple visual prompts, such as point clicks (indicating foreground/background), scribbles, or bounding boxes that roughly encapsulate the region of interest. These sparse inputs are processed by a dedicated prompt encoder—ranging from a simple MLP to a more complex Transformer encoder—to transform them into high-dimensional embedding vectors suitable for conditioning the segmentation model^{1,9}. These prompt embeddings are fused with the image features, often through cross-attention mechanisms. This fusion allows the model to take advantage of the user's global image context and local spatial information, thus providing crucial prior knowledge of the target's location and broad semantics for the segmentation task^{43,44}. This initial guidance is critical for disambiguating targets from a complex background.

Feature Fusion Stage: Refining Focus in Complex Scenarios During the feature fusion stage, visual prompts enhance the model's focus on the intended target regions, particularly when dealing with intricate or poorly defined structures in medical images. They guide the model in concentrating its representational capacity on critical areas. For instance, in the MedSAM model⁴² (as illustrated in figure 6), bounding box prompts have been shown to effectively reduce ambiguity by providing strong spatial priors, helping the model maintain segmentation accuracy even in low-contrast or structurally complex images. The

DeSAM framework⁹ employs an innovative decoupling design, combining a Prompt-Relevant IoU Module with multi-scale image features; this architecture ensures that the model can still function stably and generate accurate segmentation results even when faced with inaccurate or noisy visual prompts. ProMISE⁴³ improves a model's ability to capture 3D spatial context in volumetric medical images by introducing deep embedding layers and self-attention mechanisms specifically for processing 3D visual prompts, ensuring high efficiency and accuracy when handling complex 3D data.

Decoding Stage: Optimizing Mask Generation In the decoding stage, visual prompts continue to guide the mask decoder, ensuring it focuses on the precise target regions and thus further optimizing the generation of the final segmentation masks. After the prompt-conditioned image features are processed through the encoder and fusion stages, self-attention mechanisms within the decoder, influenced by the prompt information, direct the model to concentrate on specific regions. For example, bounding box prompts help the model to more accurately define the boundaries of target regions, which is especially critical in medical images where precise delineation impacts clinical decisions⁴². Through the coordinated application of these techniques across different stages, visual prompts not only enhance the model's initial understanding of the target in the input stage but also progressively optimize the feature fusion and decoding processes, leading to more accurate, robust, and often more efficient solutions for medical image segmentation.

Text-Guided Semantic Segmentation: Leveraging Language for Precision

Text prompts are increasingly crucial for enhancing medical image segmentation accuracy by providing rich semantic guidance. By supplying structured or natural language descriptions of anatomical structures, pathological lesions, or target characteristics, text prompts help the model better understand and precisely localize these regions of interest. This semantic conditioning typically involves several key technical stages:

Semantic Encoding of Text Prompts The process begins with converting text prompts into meaningful semantic embeddings using a text encoder. This provides strong, high-level input information to the segmentation model. In SAM-Med3D-MoE⁴⁵, for example, the encoder transforms a text prompt describing a target anatomical region into a feature vector, guiding the model to focus on that specific region. Similarly, SegVol⁴⁶ utilizes a CLIP-based text encoder to convert an input textual description (e.g., "liver") into high-dimensional vectors. These vectors are subsequently prepared for fusion with image features, laying a solid foundation for effective multimodal integration and semantic understanding.

Cross-Modal Fusion of Text and Image Features A crucial step involves effectively fusing these textual semantic embeddings with visual features extracted from the image. In the CLIP-driven Universal Model⁷ (see figure 5), text embeddings generated by CLIP, when integrated with image features (often via cross-attention), help the model to understand complex anatomical relationships (e.g., the spatial relationship between the liver and associated liver tumors). This fusion of semantic and visual information enables the model to segment specific structures or regions based on textual descriptions precisely. In LuGSAM⁴⁷, text prompts are fused with visual features extracted by a vision backbone like Grounding DINO (which itself might use a Swin Transformer), further optimizing the model's cross-modal feature integration capabilities. Through mechanisms like cross-attention, text prompts guide the model to selectively focus on specific target regions (e.g., distinguishing the right lung from the left based on textual input), thereby improving segmentation accuracy for complex or similar-looking anatomical structures.

Decoder Optimization and Output Refinement via Text Text prompts also contribute to optimizing the segmentation output during the decoder stage. In frameworks like SegVol⁴⁶, after initial fusion with image features, the text-derived information (or the fused representation) is input into the Mask decoder. This helps generate a final segmentation mask that is visually coherent and semantically aligned with the textual description, ensuring the mask accurately corresponds to the target region. In Medclip-samv2⁴⁸, text prompts are ingeniously used to generate intermediate visual prompts (such as bounding boxes or point prompts inferred from the text, often via a separate localization model or attention map analysis). These visual prompts guide an underlying segmentation model such as Segment Anything (SAM)⁴⁹, further enhancing segmentation accuracy. This indicates that text prompts can be multifaceted, as direct semantic guides and generators of more explicit spatial prompts for subsequent stages.

Spatial Localization through Hierarchical Textual prompts Beyond providing general semantic information, text prompts can significantly enhance model performance by offering precise spatial localization prompts. As demonstrated in Ariadne's thread⁵⁰, text prompts can be structured to gradually provide location information about the target region hierarchically, progressing from coarse to fine descriptions. This might include broad regional descriptions (e.g., "left lung"), more specific location descriptions (e.g., "upper lobe infection of the left lung"), and finally, finely detailed lesion characteristics (e.g., "infection area within the superior segment of the upper lobe"). These multi-level spatial prompts, embedded in text, help the model more accurately identify and localize the target region, thereby addressing precision issues in image segmentation that often arise from insufficient contextual information or ambiguity.

Facilitating Weakly Supervised Learning and Pseudo-Label Generation Text prompts play a vital role in weakly supervised learning paradigms and the generation of pseudo-labels, reducing the reliance on extensively manually annotated data. In Simtxtseg⁵¹, text prompts are used to help generate initial visual prompts or attention maps, which are then further processed to create pseudo-masks. These pseudo-masks serve as training targets for the segmentation network, which is particularly valuable in scenarios where pixel-level annotations are scarce. Similarly, TP-DRSeg⁵² employs an explicit prior encoder to transform text descriptions of diabetic retinopathy lesions into visual priors. These priors are then fused with image features to enhance segmentation accuracy. This approach, which improves the model's generalization ability through pseudo-label generation and weakly supervised training, further underscores the multifaceted utility of text prompts in modern segmentation networks, especially for specialized medical tasks.

Advanced Prompting Strategies: Learnable, Adaptive, and Unsupervised Approaches

Beyond direct visual and textual inputs, medical image segmentation has witnessed the development of more sophisticated prompting strategies. These often involve learnable components, adaptive mechanisms, or approaches tailored for unsupervised or domain-adaptation challenges, aiming further to boost segmentation accuracy, robustness, and applicability.

Learnable Prompts for Task-Specific Adaptation The core idea behind "learnable prompts" is to obtain a set of parameterized vectors through end-to-end training. These vectors implicitly encode task- and modality-specific prior knowledge, guiding network behavior more effectively than fixed prompts for specific segmentation tasks. In⁴, learnable prompt vectors are injected into UNet-style "promptable blocks," cascaded with windowed image embeddings, and jointly process image features and task information via multi-head attention mechanisms. The Medunise model⁵³ separately employs modality prompts (injected at early encoder stages to help differentiate multimodal input data) and task prompts (fused with sample features via a FUSE module at the encoder's terminal to provide task-specific priors for the decoder). Uniseg¹⁴ focuses on capturing intrinsic correlations between different segmentation tasks by integrating universal prompts with encoder outputs, delivering task-specific priors at the initial decoder stage. The DCTP-Net⁵⁴ introduces a Learnable Prompt block within its prompt-aware branch specifically to extract brain prior knowledge from CT images, which is then fused with image features to reduce structural detail interference during acute ischemic stroke lesion segmentation. The FedLPPA framework⁵⁵ for federated learning constructs a "triple prompt" system comprising globally shared universal knowledge prompts, locally distributed data-aware prompts, and predefined annotation-sparsity prompts. These diverse prompts are fused with encoder outputs through a Tri-prompt Dual-attention Fusion module and injected with rich contextual information before decoding, enabling personalized and efficient federated weakly-supervised segmentation.

Prompts for Unsupervised Learning and Domain Adaptation Several studies have introduced novel prompting strategies to tackle unsupervised annotation or domain adaptation challenges in medical image segmentation. Chen et al.⁸ propose the Self-prompt mechanism, which employs a multi-scale self-prompt generation module. By fusing features from various layers of a domain-adaptive encoder with a Feature Pyramid Network for foreground prediction and subsequent screening, high-confidence (self-generated) prompts are element-wise added to the final image features. These are fed into a domain query-enhanced decoder, achieving precise segmentation of nuclear regions without manual annotation. Lin et al.⁵⁶ introduce a Domain-specific prompt embedded within a feature transfer module. This module automatically extracts and fuses unique information pertinent to the current target domain, guiding the network to generate domain-aware features and enabling robust feature alignment and adaptation across multiple target domains. Na et al.⁵⁷ present the Auto-prompt system, which utilizes an independent auxiliary network to generate initial prediction masks. After Sigmoid activation, high-confidence positive and negative prompt points are automatically extracted from these masks and fused with SAM's Prompt Encoder. This guides the Mask Decoder in completing the cell nucleus segmentation task more accurately, automating the prompt generation process. Luo et al.⁵⁸ propose the Task-Specific Prompt as a trainable prompt vector embedded at the back end of a Vision Transformer (ViT) based segmentation model. Together with the image features extracted by the encoder, this prompt enters the Transformer decoder to provide task-specific prior information, tailoring the generic ViT to specific segmentation needs. Meanwhile, Chen et al.⁵⁹ introduce the Low-frequency prompt. This technique modulates the low-frequency amplitude of test images during the preprocessing stage. The goal is to adjust the image style and texture in the frequency domain, making the test images appear closer to the characteristics of the source domain on which a pre-trained, frozen segmentation network was trained, thus mitigating the detrimental impact of domain shift.

Prompt Mechanisms in Medical Image Classification

Prompt-based methodologies have become essential for enhancing the performance, generalization capabilities, and interpretability of models engaged in medical image classification. By effectively leveraging prompts, these models can better navigate the inherent complexities of medical data, thereby addressing critical challenges such as domain adaptation, the scarcity of labeled data (enabling few-shot or zero-shot learning), and the increasing demand for explainable AI in clinical decision support. This section details the primary categories of prompts employed for these classification tasks.

Text-Driven Classification: From Static Labels to Dynamic Semantic Guidance

Text prompts integrated within medical image classification models, predominantly processed by the text encoder component of vision-language architectures, have evolved significantly. They have transitioned from simple static class descriptors to dynamic, trainable, and knowledge-rich constructs, substantially enhancing domain specificity, model interpretability, and overall classification accuracy, particularly in zero-shot and few-shot learning paradigms.

Foundational Approaches with Static and Pre-trained Prompts Early and foundational frameworks, such as that detailed in¹¹, utilize static text prompts, typically defined by class names (e.g., “well differentiated tubular adenocarcinoma”) sourced directly from dataset labels. These prompts are fed into a frozen biomedical language model (e.g., BioLinkBERT), which generates fixed class embeddings x_{cl} . These embeddings remain invariant during training (due to a “stop-grad” operation) and serve as stable semantic reference points. They are then compared (often via cosine similarity) with image embeddings x'_i , which are derived from a vision encoder (e.g., ViT-B/16) and projected into a shared latent space. The model is optimized using a cross-entropy loss to align these visual features with their corresponding textual semantics, facilitating classification without altering the pre-trained weights of the text encoder.

This static prompting approach is further refined and scaled in influential models like BiomedCLIP⁶⁰ and PubMedCLIP⁶¹. In these contexts, text prompts often consist of richer textual data, such as PubMed captions (e.g., “Histology of metastatic amelanotic melanoma...”²¹) or excerpts from medical reports (e.g., “CT of the chest showing multiple thick-walled cavities...”⁶²). These are processed by specialized medical text encoders (e.g., PubMedBERT²³, BioClinicalBERT²²). During large-scale pre-training, these prompts drive contrastive learning (e.g., using InfoNCE loss) to align the embeddings from the text encoder with those from the vision encoder by iteratively optimizing their similarity scores²¹. For downstream zero-shot classification tasks, the pre-trained text encoder converts new class name prompts into target embeddings t_p . These are then matched against the vision encoder’s output v_p for a given image using cosine similarity, enabling classification inference on unseen classes without task-specific retraining²². The text encoder is a crucial conduit for injecting domain-specific prior knowledge. In contrast, the vision encoder adapts to extract relevant visual features, a synergy often harmonized by a projection head that aligns the dimensions of the different modalities.

Advanced Dynamic and Learnable Text Prompts for Enhanced Adaptability Moving beyond static inputs, more recent frameworks like XCoOp⁶³ and MSCPT⁶⁴ introduce dynamic and tunable (learnable) text prompts, significantly increasing their operational complexity and adaptability within the text encoder. In XCoOp, text prompts comprise a combination of soft prompts (learnable vectors, often initialized with a template like “a photo of a [disease name]”) and clinical challenging prompts (fixed textual descriptions derived from experts, for example, “a photo of melanoma, with an irregular pigment network...”), both of which are processed by the text encoder. The soft prompts are optimized at both token and prompt levels, aligning their learned embeddings V with fixed clinical embeddings Q (derived from challenging prompts) using contrastive and cross-entropy losses. This process enhances interpretability by explicitly tying each learnable token to a specific clinical concept. These rich prompt embeddings then interact with the vision encoder’s outputs—both global image features p and local patch features $F = \{f_1, f_2, \dots, f_M\}$ —through a global-local alignment loss, refining the classification process by mimicking expert diagnostic reasoning which often involves both holistic assessment and attention to local details. MSCPT⁶⁴ further leverages a dual-path text encoder system for multi-scale analysis of Whole Slide Images (WSIs). Frozen low-level encoders process prompts relevant to low-magnification views (e.g., 5×, with prompts like “Sparse stroma among tumor cells”), while learnable (prompted) high-level encoders handle prompts for high-magnification views (e.g., 20×, with prompts like “Distinct cellular borders...”). The resulting multi-scale textual embeddings are then integrated using transformer layers. This hierarchical textual processing feeds into graph-based and non-parametric pooling mechanisms, aligning with vision encoder outputs to capture contextual and fine-grained multi-scale features essential for accurate WSI classification.

Text Prompts in Multiple Instance Learning Paradigms In the context of Multiple Instance Learning (MIL)^{65,66}, which is prevalent in computational pathology where an entire WSI (a “bag”) is classified based on its constituent patches (“instances”), text prompts extend their influence beyond simple text encoding to orchestrate both instance-level feature extraction and bag-level aggregation.⁶⁵ employs composite prompts: these include task labels (e.g., “a WSI of [Lung adenocarcinoma]”), detailed descriptions generated by models like GPT-4 (e.g., “glandular or acinar formations”), and learnable vectors. The text encoder processes all these textual components to generate instance prototypes P and bag-level tokens B_i . These textual embeddings then guide the aggregation of patch features Z_i (from the vision encoder) via MIL pooling mechanisms, culminating in a bag-level similarity matching for final classification. Similarly,⁶⁶ optimizes textual prompts (e.g., “An H&E image of [breast adenocarcinoma]”) within the text encoder to produce embeddings \hat{z}_t . These embeddings \hat{z}_t are then used to weight the vision encoder outputs \hat{z}_v from individual patches in a context-driven pooling scheme, often enhanced by a residual visual adapter to capture visual nuances better. In both these MIL frameworks, the embeddings derived from text prompts dynamically steer the vision encoder’s feature extraction and, more critically, the aggregation process, adapting effectively to low-shot learning scenarios and improving the classification of complex, heterogeneous medical images like WSIs.

Specialized Frameworks Balancing Adaptability and Simplicity The operational depth and utility of text prompts are further exemplified in dedicated frameworks like FLAIR⁶⁷ and MI-Zero⁶⁸, which strive to balance model adaptability with implementational simplicity for specific medical imaging domains. FLAIR utilizes expert-derived, descriptive textual prompts (e.g., “only a few microaneurysms are present” for retinal images) encoded by a BioClinicalBERT text encoder. These embeddings are used for contrastive pre-training and subsequent zero-shot inference, aligning with vision encoder outputs via similarity scores to perform classification¹³. Conversely, MI-Zero employs static class-name templates (e.g., “an image showing adenocarcinoma”) within a standard text encoder like PubMedBERT²³. This produces fixed embeddings w_m used for patch-level similarity computation, aggregating results via MIL operators such as topK pooling for WSI classification⁶⁸. While FLAIR’s approach allows for tunable knowledge integration through rich, expert-crafted prompts, MI-Zero’s strength lies in its simplicity and reliance on robust pre-trained alignments, with both methods effectively interfacing with vision encoders to drive classification.

Diverse Non-Textual and Structurally-Informed Prompts for Classification

Beyond text-based semantic guidance, various innovative non-textual and structurally-informed prompt mechanisms have been developed to address specific challenges in medical image classification. These often target domain generalization, memory efficiency, leveraging spatial or structural image information, or enhancing model reasoning capabilities.

Learnable Prompts for Domain Generalization and Efficiency To tackle the critical problem of domain generalization,⁵ proposes Collaborative Domain Prompts. These are lightweight, learnable vectors integrated into the input layer of ViT, typically prepended to the CLS token and patch embeddings. By leveraging a shared prompt component and a domain-specific prompt generator, these prompts facilitate collaboration and knowledge sharing between representations learned for different domains, enabling robust generalization to unseen domains in medical image classification, a crucial capability when explicit domain labels are unavailable or undefined. Addressing memory efficiency in incremental learning scenarios,¹² introduces Domain Prompts for histopathology classification. Applied within ViT’s multi-head self-attention layers, these prompts are designed to decouple into domain-specific and domain-invariant components. This allows the model to capture unique domain features and shared knowledge across incrementally learned tasks with minimal memory overhead, preventing catastrophic forgetting.

Structural and Embedding-Based Prompts for Fine-Grained Analysis For tasks requiring fine-grained analysis, such as nucleus detection and classification in pathology,⁶⁹ presents nucleus group embeddings. These are learnable embeddings added to the input space of Swin Transformers, serving as “grouping prompts.” They provide initial representations within a grouping transformer classifier to capture semantic similarities among nuclei, enhancing detection and subsequent classification accuracy.

Zhu et al.⁷⁰ utilize Segmentation Map Prompts for classification tasks where spatial context is important. A dedicated encoder encodes pre-computed or predicted segmentation maps (e.g., identifying kidney stones) into token representations. These tokens are then fed alongside the image tokens into a frozen ViT. The segmentation map prompts interact with image features through self-attention mechanisms, guiding the model to focus on relevant regions delineated by the map, thereby improving classification accuracy by explicitly incorporating spatial context. Further improving adaptability in nucleus classification,⁷¹ integrates Label Prompts within a dynamic prompt module between the feature encoder and the decoder. These prompts, derived from class labels, help to refine multi-scale features, improving the adaptability and accuracy of nucleus classification across diverse datasets.

Generative, Prototypical, and Inferential Prompts for Advanced Classification Leveraging generative capabilities,⁷² proposes Pseudo-Prompts, which constitute dynamically generated class-specific embeddings produced by a dedicated prompt decoder module. These pseudo-prompts are subsequently fed into a frozen text encoder to facilitate multi-label zero-shot learning. The frozen text encoder operates without learnable embedding layers for these specific inputs, effectively leveraging the priors learned by pre-trained vision-language models. This approach enables the system to harness established knowledge representations while maintaining the flexibility to generate task-specific prompt embeddings for novel classification scenarios.

To better capture diverse pathological characteristics in WSIs,⁷³ develops Prototypical Visual Prompts. These are cluster-center-based prompts derived from representative image patches identified before feature extraction. By conditioning the model on these visual prototypes, they enrich WSI classification. Enhancing model reasoning,⁷⁴ combines Few-shot Learning with Chain of Thought Prompts at the input level of large vision-language models (e.g., GPT4) for tasks like blood cell classification. This approach aims to improve classification through enhanced contextual learning and inferential capabilities by providing examples and explicit reasoning steps. However, performance can vary significantly depending on the base model’s capacity.

THE APPLICATION OF THE PROMPT MECHANISM

Prompt engineering enables large vision models to achieve a wide range of AI tasks in the medical field with high performance, as shown in figure 7. With customized medical prompts, models can solve complex medical image problems that were

previously difficult without giant medical data sets and resources. By designing appropriate prompts, developers can quickly obtain target data without training or fine-tuning the model, which is time-consuming and costly. In the foreseeable future, prompting engineering (especially text prompts) will be further developed to utilize the potential of large medical imaging models fully. In the medical image field, the applications of prompt engineering include:

- **Generation:** Introducing prompt engineering into large-scale medical image models offers boundless potential for generating images from textual or report prompts. These advancements hold promise for addressing the challenge of limited medical imaging data in deep learning and applications in medical education and training.
- **Segmentation:** Medical imaging datasets often suffer from small size and partially labeled problems. However, by leveraging prompts, models can be trained on large amounts of diverse medical data and learn the associations between data, such as anatomical structures, enabling the development of powerful and universal medical segmentation models.
- **Classification:** In medical image classification tasks, prompts are commonly employed to address the issue of sample scarcity, such as in zero-shot and few-shot learning. Furthermore, prompts, especially learning-based prompts, are employed to enhance the classification performance of models.

We also provided a more comprehensive overview of applications of prompt engineering in medical images, as summarized in Table 1, Table 2, Table 3, Table 4, and Table 5.

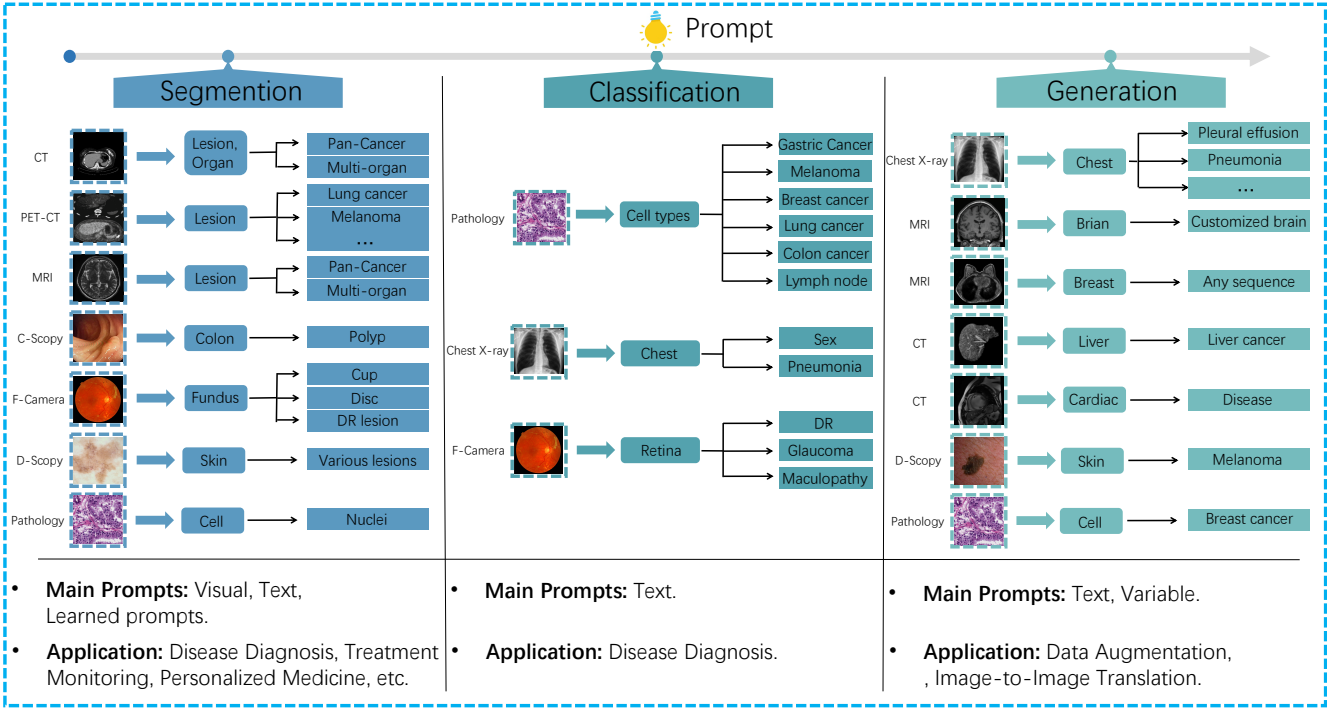


Figure 7. Typical applications of prompts in three major medical imaging tasks: segmentation, classification, and generation. Note: P represents positive cases, N represents negative cases, C-scopy refers to colonoscopy, DR refers to Diabetic Retinopathy and D-scope refers to dermatoscope.

Medical Image Generation

One of the significant challenges in developing large vision models in the medical imaging field is the lack of high-quality training data. This is often attributed to the challenges in accessing medical data, which stem from its scarcity and the privacy concerns associated with medical imaging data. In addition, constructing large-scale, accurately annotated datasets requires much effort from experienced radiologists, which is usually impossible. Medical imaging data sets are often imbalanced as pathologic findings are generally rare, which also hurts the training effect of the model¹⁶³. The combination of Large-scale Model (Foundation Model) and prompt engineering offers promising prospects for addressing this issue. In addition to developing expert models, the generated image data can be used for medical education and training. Text-prompt image generation could generate realistic training scenarios for medical students and healthcare professionals. By simulating a diverse

Table 1. Applications of Prompt Engineering in Generation

Reference	Year	Task	Prompt Type	Link
Akrout et al. ²	2023	Skin	Text	-
Pinaya et al. ³	2022	Brain MRI	Multivariable (age, sex, and brain information)	-
Chambon et al. ⁶	2022	Chest X-ray	Text	-
Xu et al. ²⁰	2024	Chest CT	Text	-
Yellapragada et al. ¹⁹	2024	Histopathology	Text	https://github.com/cvlab-stonybrook/PathLDM
Moghadam et al. ⁴¹	2023	Histopathology	Label	-
Hamamci et al. ⁷⁵	2025	Chest CT	Text	https://github.com/ibrahimethemhamamci/GenerateCT
Dai et al. ²⁹	2024	Joint CT	Text	-
Wang et al. ⁷⁶	2024	Brain MRI	MR imaging parameters	-
Shi et al. ⁷⁷	2025	Breast Ultrasound	Texture structure and lesion information	-
Dahan et al. ⁷⁸	2024	Musculoskeletal Ultrasound	Image	-
Bluthgen et al. ³⁰	2024	Chest X-ray	Text	-
Yu et al. ⁷⁹	2024	Abdominal Lymph Node CT	Mask	-
Chambon et al. ³¹	2022	Chest X-ray	Text	-
Xiao et al. ⁸⁰	2023	Liver CT	Mask	-
Weber et al. ⁸¹	2023	Chest X-ray	Text	-
Hashmi et al. ³²	2024	Chest X-ray	Text	https://github.com/BioMedIA-MBZUAI/XReal
Liang et al. ³³	2024	Chest X-ray	Text	-
Han et al. ³⁴	2024	Chest X-ray	Text	https://micv-yonseio.github.io/cxrl2024/
Shentu et al. ⁸²	2024	Chest X-Ray	Combining text embedding and image embedding	https://github.com/junjie-shentu/CXR-IRGen
Liu et al. ³⁵	2024	4D Cardiac Cine MRI	Text	https://github.com/me-congliu/TexDC
Borghesi et al. ⁸³	2024	Skin	Mask	https://github.com/aequitas-aod/experiment-gen-skin-images
Fang et al. ⁸⁴	2024	Chest X-Ray	Edge maps, depth maps, Masks, or CLIP image embeddings	-
Sagers et al. ⁸⁵	2022	Skin	Text	-
Wang et al. ⁸⁶	2025	Brain MRI	Text	https://github.com/Wangyulin-user/TUMSyn
Li et al. ⁸⁷	2025	Brain MRI	Localization Prompt	https://github.com/ChanghuiSu/TLP
Duan et al. ⁸⁸	2025	Fetal Ultrasound	Anatomical Prompt	https://dyf1023.github.io/FetalFlex/
Liu et al. ⁸⁹	2025	Brain MRI	multi-parametric MRI and treatment information	-

Table 2. Applications of Prompt Engineering in Segmentation (Part 1)

Reference	Year	Task	Prompt Type	Link
Mattjie et al. ¹	2023	Skin Lesions, Lung, Polyps, Breast Tumor, Femur and Ilium	Visual	https://github.com/Malta-Lab/SAM-zero-shot-in-Medical-Imaging
Fischer et al. ⁴	2024	Pancreas-CT, Cranial Vault abdomen CT	Learnable prompt	https://github.com/marcdcfischer/PUNETR
Liu et al. ⁷	2023	Organ and Tumor Detection	Text	https://github.com/ljwztc/CLIP-Driven-Universal-Model
Chen et al. ⁸	2024	Nuclei	Self-prompt	https://github.com/CUHK-AIM-Group/UN-SAM
Gao et al. ⁹	2023	Prostate	Visual	https://github.com/yifangao112/DeSAM
Cheng et al. ¹⁰	2023	Universal Segmentation	Visual, Masks	https://github.com/uni-medical/SAM-Med2D
Ye et al. ⁵³	2024	Universal Segmentation	Learnable prompt	https://github.com/yeerwen/UniSeg
Wu et al. ⁹⁰	2024	Organ	Visual, Mask	https://github.com/KidsWithTokens/one-prompt
Li et al. ⁴³	2024	CT pancreas and colon	Visual	https://github.com/MedICL-VU/ProMISe
Lin et al. ⁵⁶	2023	Infant MRI of different ages, high- and low-grade gliomas	Domain-specific prompt	https://github.com/MurasakiLin/prompt-DA
Chang et al. ⁹¹	2023	Organ	Visual	-
Bai et al. ⁹²	2023	CT liver tumor	Visual	-
Wu et al. ⁹³	2024	Dermatology image, optical disc, Chest X-ray	Image	-
Ma et al. ⁴²	2024	Universal Segmentation	Visual	https://github.com/bowang-lab/MedSAM
Xu et al. ⁹⁴	2023	Universal Segmentation	Pixel Similarities	-
Xie et al. ⁹⁵	2024	CT organ	Visual, Mask	-
Ramesh et al. ⁴⁷	2024	Chest X-Rays	Text, Visual	-
Kato et al. ⁹⁶	2023	Cell	Visual	https://github.com/usagisukisuki/oneshot-part-cellsegmentation
Zhang et al. ⁹⁷	2023	Organ and Tumor Detection	Text	https://github.com/MrGiovanni/ContinualLearning
Ye et al. ¹⁴	2023	Organ and Tumor	Learnable Prompt	https://github.com/yeerwen/UniSeg
Deng et al. ⁴⁴	2023	Optic Cup and Disc in Fundus Photographs	Visual	-
Na et al. ⁵⁷	2024	Nuclei	Auto-prompt	-
Zhong et al. ⁵⁰	2023	Chest X-Rays	Text	https://github.com/Junelin2333/LanGuideMedSeg-MICCAI2023
Du et al. ⁴⁶	2023	Organ	Visual, Text	https://github.com/BAAI-DCAI/SegVol
Tomar et al. ⁹⁸	2022	Organ	Visual, Mask	https://github.com/nikhilroxtomar/TGANet
Xie et al. ⁵¹	2024	Polyp, MRI Brain	Text, Visual	-
Koleilat et al. ⁴⁸	2024	Universal Segmentation	Text, Visual	https://github.com/HealthX-Lab/MedCLIP-SAM
Zhao et al. ⁹⁹	2023	Universal Segmentation	Text	https://github.com/zhaoziheng/SAT

Table 3. Applications of Prompt Engineering in Segmentation (Part 2)

Reference	Year	Task	Prompt Type	Link
Biswas et al. ¹⁰⁰	2023	Polyp	Visual, Mask	https://github.com/RisabBiswas/Polyp-SAM++
Chen et al. ¹⁰¹	2025	Abdominal and cardiac MRI to CT	Visual	-
Han et al. ¹⁰²	2023	Nuclear, chest X-ray, Colon histology	Visual	-
Saeed et al. ¹⁰³	2023	PET-CT Head and Neck Cancer	Learnable parameters, Termed prompts	-
Li et al. ¹⁰⁴	2023	CT scans liver and pancreatic lesions	Visual	-
Luo et al. ⁵⁸	2023	Universal Segmentation	Task-Specific Prompt	-
Xie et al. ¹⁰⁵	2023	Abdominal CT and cardiac MRI	Visual	-
Chen et al. ¹⁰⁶	2023	Nuclei	Visual	https://github.com/xq141839/SPPNet
Sridhar et al. ¹⁰⁷	2023	Chest X-ray	Mask, Visual	-
Zhang et al. ¹⁰⁸	2023	CT pulmonary nodule	Visual	-
Glatt et al. ¹⁰⁹	2023	Liver cell	Visual	-
Zhou et al. ¹¹⁰	2024	Pancreas-CT, Brain MRI	Cross-Prompt	https://github.com/Fyw1988/MUL SCIP
Huang et al. ¹¹¹	2024	Universal Segmentation	Visual prompt	-
Chen et al. ⁵⁹	2024	Joint optic disc (OD) and cup (OC), polyp	Low-frequency prompt	https://github.com/Chen-Ziyang/VPTTA
Ouyang et al. ¹¹²	2024	Universal Segmentation	Text prompt	-
Chen et al. ¹¹³	2024	MR femur, tibia, femoral- and tibial cartilages	Visual	https://github.com/chrissyinreallife/KneeSegmentWithSAM.git
Shaharabany et al. ¹¹⁴	2024	CT organ, nuclear	Mask, Visual	https://github.com/talshaharabany/ZeroShotSAM
Wang et al. ¹¹⁵	2024	MRI colon and stomach, CT Liver Tumor	Temporal prompt	-
Adhikari et al. ¹¹⁶	2024	Universal Segmentation	Visual	-
Kong et al. ¹¹⁷	2024	Universal Segmentation	Visual	https://github.com/naamiinepal/tunevlseg
Liu et al. ¹¹⁸	2024	Kidney pathology	Automatic prompt	https://github.com/SnowRain510/GBMSeg
Lin et al. ⁵⁵	2024	Fundus OD/OC, OCTA FAZ, Endoscopy Polyp, MRI Prostate	Learnable prompt	https://github.com/llmir/FedLPPA
Chen et al. ¹¹⁹	2024	Universal Ultrasound	Task Prompt	-
Cui et al. ¹²⁰	2024	Kidney pathology	Text, Visual	-
Xie et al. ¹²¹	2024	Polyp	Visual	-
Xia et al. ¹²²	2024	MRI Cervical cancer	Visual	-
Zhu et al. ⁷⁰	2024	Kidney stone	Mask	-
Teng et al. ¹²³	2024	MRI brain	Text	https://github.com/TL9792/KGPL
Chen et al. ¹²⁴	2024	Cine Cardiac MRI	Visual	-
Guan et al. ¹²⁵	2024	Universal Segmentation	Class and Visual	-
Song et al. ¹²⁶	2024	Histopathological breast cancer	Mask	https://github.com/QI-NemoSong/EP-SAM
Khor et al. ¹²⁷	2024	CT Nasopharyngeal Carcinoma with Prior Anatomical Information	Visual	-

Table 4. Applications of Prompt Engineering in Segmentation (Part 3)

Reference	Year	Task	Prompt Type	Link
Xue et al. ¹²⁸	2024	MRI cervical cancer	Visual	-
Cui et al. ¹²⁹	2024	Nuclei	Visual	-
Lyu et al. ¹³⁰	2024	CT liver tumor	Visual	-
Yang et al. ¹³¹	2024	Chest X-ray	Task-adaptive Visual	-
Xue et al. ¹³²	2024	CT Endometrial cancers	Visual	-
Dai et al. ¹³³	2024	CBCT dental images	Visual	-
Cui et al. ¹³⁴	2024	Kidney pathology	Text	-
Hu et al. ¹³⁵	2024	Polyp, colonoscopy, skin injury	Visual	-
Liu et al. ⁵⁴	2024	CT Acute Ischemic Stroke Lesion	Learnable prompt	-
Sun et al. ¹³⁶	2024	MRI Brain Tumor	Tumor grades	https://github.com/YonghengSun1997/AEPL
Song et al. ¹³⁷	2024	Laryngoscopic Image	Visual	https://github.com/youcongzh
Cheng et al. ¹³⁸	2024	MRI brain	Frequency filtering and spatial perturbation prompts	-
Li et al. ¹³⁹	2024	Nucleus	Visual	-
Zhang et al. ¹⁴⁰	2024	Multi-Class Cell	Text	-
Huang et al. ¹⁴¹	2024	Tooth	Image	-
Li et al. ¹⁴²	2024	MRI Breast Cancer Tumor	Learnable prompt	-
Shan et al. ¹⁴³	2025	Polyp, Chest X-ray, Chest CT	Text prompt	https://github.com/HUANGLIZI/STPNet
Wang et al. ¹⁴⁴	2025	Brain tumor MRI, Abdominal CT, Cardiac MRI	Mask prompt	https://github.com/wanghr64/WeakMedSAM
Yin et al. ¹⁴⁵	2025	breast ultrasound	Automatic prompt	-
Liu et al. ¹⁴⁶	2025	Joint optic disc (OD) and cup (OC), Brain tumor MRI, Hole heart CT and MRI	Deformable Convolutional Prompt	-
Yin et al. ¹⁴⁷	2025	Multi-organ abdominal CT and MRI	Frequency prompt	-
Gao et al. ¹⁴⁸	2025	Total-Body PET	Textual and disentangled organ features	-
Tian et al. ¹⁴⁹	2025	Plaque and vessel carotid ultrasound	Self-prompt	-
Zou et al. ¹⁵⁰	2025	Prostate MRI	Advanced Prompt Points	-
Chen et al. ¹⁵¹	2025	Nuclei	Self-prompt	https://github.com/CUHK-AIM-Group/UN-SAM
Zhang et al. ¹⁵²	2025	Nuclei	Category Prompt	https://github.com/CUHK-AIM-Group/UN-SAM
Zhao et al. ¹⁵³	2025	Skin Lesion, Thyroid Ultrasound, Spine CT, Cardiac MR	Auto edge prompt	-

Table 5. Applications of Prompt Engineering in Classification

Reference	Year	Task	Prompt Type	Link
Yan et al. ⁵	2024	Melanoma classification and Cancerous tissue detection, Diabetic Retinopathy classification	Collaborative domain prompts	https://github.com/SiyuanYan1/PLDG/tree/main
Zhang et al. ¹¹	2023	Pathological image	Text, Learning Visual Prompt	https://github.com/Yunkun-Zhang/CITE
Zhu et al. ¹²	2024	Breast cancer, epithelium-stroma tissue histopathology	Domain prompt	-
Huang et al. ⁶⁹	2023	Nucleus Detection and Classification	Nucleus group embeddings	-
Qu et al. ⁶⁵	2024	Breast cancer, lung cancer, and cervical cancer Pathology	Text	https://github.com/miccaiif/TOP
Guo et al. ¹⁵⁴	2023	Lesion	Multiple text	-
Wang et al. ²²	2022	Disease Classification	Text	https://github.com/RyanWangZf/MedCLIP
Lu et al. ⁶⁸	2023	Histopathology	Text	https://github.com/mahmoodlab/MI-Zero
Zhu et al. ⁷⁰	2024	Kidney Stone	Segmentation Map	-
Huang et al. ⁷¹	2024	Nucleus Classification	Label prompt	https://github.com/lhaof/UniCell
Eslami et al. ⁶²	2021	Disease type, tumor location, etc.	Text	https://github.com/sarahESL/PubMedCLIP
Zhang et al. ²¹	2023	Lung tissue, colon tissue, pneumonia	Text	-
Silva et al. ⁶⁷	2023	Retina	Text	https://github.com/jusiro/FLAIR
Cao et al. ¹⁵⁵	2023	Brain Tumor, Chest disease	Learnable contexts	-
Zheng et al. ¹⁵⁶	2024	Chest X-ray disease	Automatic prompt	-
Ye et al. ⁷²	2024	Chest X-ray	Pseudo-Prompt	https://github.com/fallingnight/PSPG
Huang et al. ¹⁵⁷	2024	Fundus, dermoscopic, mammography, chest X-ray	Fine-grained prompt from pre-trained models	-
Lin et al. ⁷³	2024	Breast cancer histopathology	Prototypical Visual Prompt	-
Chikontwe et al. ⁶⁶	2024	Colorectal cancer, breast cancer metastasis detection in lymph nodes histopathology	Text	-
Han et al. ⁶⁴	2024	Lung , breast and Kidney cancer histopathology	Text	-
Yang et al. ¹⁵⁸	2024	Gastric adenocarcinoma histopathology	Text-augmented Visual Prompt	-
Sanchez et al. ⁷⁴	2024	Blood cell	Few-shot learning, Chain of thought	-
Bai et al. ¹⁵⁹	2025	General	Label-Semantic-Based Prompt	-
Koleilat et al. ¹⁶⁰	2025	General	Learnable prompt	https://github.com/HealthX-Lab/BiomedCoOp
He et al. ¹⁶¹	2025	General	Dynamic Visual Prompt	-
Luo et al. ¹⁶²	2025	Skin, Gastrointestinal and Respiratory disease	Decoupled Probabilistic Prompt	https://github.com/CUHK-AIM-Group/LDPP

array of cases, rare conditions, and complex anatomical variations, synthetic medical image generation could enrich educational programs and offer invaluable hands-on experience in a controlled environment¹⁶⁴.

Expanding on the growing success of Stable Diffusion¹⁶⁵ and GAN¹⁶⁶ in image generation, an increasing number of researchers are integrating these models into the domain of medical image generation to address the inadequacies stemming from limited medical imaging data. While previously, the generating ability of diffusion models was mostly used for unconditional generation of data, such as^{167–169}, these models usually only have the ability to solve a single problem and have low flexibility, making it difficult to transfer to other tasks. The proposal of Foundation Models such as CLIP and DALL-E has greatly promoted the development of text-to-image and image-to-text fields. Benefiting from the development of these multi-modal models, more and more models are being developed specifically for medical image generation, which will help develop more efficient tools to serve radiologists and patients.

Text prompts play a crucial role in image generation by offering flexibility, interpretability, and the ability to define specific disease types, anatomical structures, or imaging features. This enables the generation of highly customized medical images tailored to particular needs. A representative task in this domain is generating Chest X-rays using medical reports^{6,30–34,81,82,84}, with an illustrative figure adapted from³⁴, as shown in figure 4. Beyond this, various tasks have also been explored, including skin image generation^{2,83}, MRI generation^{3,35,76}, histopathology image generation^{19,41,85}, and CT image generation^{20,75,79,80}. Apart from text prompts, depending on the specific application, prompts can take various forms such as conditional variables³, masks⁸⁰, and labels⁴¹.

In summary, broad prompts include text and variables, sequences, classes, etc, used to control model output. For example, in³, age, sex, and brain structure volumes are utilized as prompts to generate the expected brain images.⁴¹ generating histopathology images with a genotype prompt. However, text prompts remain the primary focus for future development due to their higher scalability and lower usage threshold.

Medical Image Segmentation

Accurate medical image segmentation can improve diagnostic accuracy, treatment planning, and disease monitoring⁴². In image segmentation, textual and visual prompts are among the most widely applied approaches. Textual prompts are commonly transformed into encoded embeddings by language models and then fed into segmentation models. Visual prompts, including points, boxes, and masks, are typically utilized in Visually Prompted Models. A representative work in this area is SAM, which has quickly triggered the development of large models for medical image segmentation.

The introduction of CLIP¹⁷⁰ has highlighted the enormous potential of textual prompts, reshaping how researchers approach multimodal learning and image segmentation tasks. Among the representative works,⁷ stands out as a landmark study that developed a CLIP-driven universal model for organ and tumor segmentation, as shown in figure 5. In this work, CLIP was utilized to generate segmentation prompts, which, compared to traditional one-hot prompts¹⁷¹, capture anatomical relationships more effectively and significantly expand the dataset. Building on this work,⁹⁷ further optimized the model by assigning a separate, independent MLP to each class, which reduced interference between different classes. However, CLIP has limited ability to generalize in medical scenarios due to the differences between natural and medical texts¹⁴. Some studies^{50,99,172} employ contrastive learning or utilize specialized text encoders^{23,24,173} specifically designed for medical images.

SAM, as a representative large-scale model in the field of image segmentation, has inspired many studies^{1,8,9,44,57,94} to focus on enhancing visual prompts for directly segmenting medical images using the SAM model. These methods have significantly improved SAM's capabilities in medical image segmentation. Although SAM demonstrates strong segmentation quality and zero-shot generalization to novel scenes and unseen objects, its training data does not include medical images. Consequently, its performance in most medical image segmentation tasks is often unsatisfactory^{174–177}. Therefore, many researchers have focused on developing large models specifically for medical image segmentation. For example,^{42,178} collected large volumes of medical image data and fine-tuned it based on SAM, resulting in a large-scale promotable model designed for medical image segmentation. The MedSam was shown in figure 6. It also explored the effects of different prompts on segmentation performance. Box-based prompts provided more explicit guidance, whereas point-based prompts were more prone to ambiguity in medical image segmentation. In¹⁷⁹, the SAM-Med2D model was extended to SAM-Med3D, where sparse prompts were enhanced with 3D position encodings to capture spatial nuances, while dense prompts were processed using 3D convolutions.

There is also work combining text and visual prompts. For instance,⁴⁷ extracted text features from text prompts using BERT¹⁸⁰, and then created bounding boxes that were used as prompts for SAM. Similarly,⁴⁸ proposed a novel framework called MedCLIP-SAM, which combines CLIP and SAM models to generate segmentations of clinical scans. They fine-tuned the BiomedCLIP model, and the visual prompts were generated using image and text, post-processed with gScoreCAM¹⁸¹. Other similar approaches, such as^{51,52,100}, leverage textual prompts to generate visual prompts that enhance segmentation tasks in medical imaging. Some works, like⁴⁶, directly fuse textual and visual prompts to improve segmentation. In addition to visual and textual prompts, other types of prompts have been explored. For example,⁴ proposed a promptable UNet architecture that

adapts to segmentation tasks using class-dependent, learnable prompt tokens. Similarly,^{14,53} designed a learnable universal prompt that captures the relationships among tasks and converts it into task-specific prompts, which are fed into the decoder as part of its input.

Visual prompts are often applied to large foundation models such as SAM and its derivatives. These prompts can directly indicate areas of interest in the image, providing a more intuitive and precise approach than textual descriptions. Text prompts, on the other hand, can sometimes lead to inconsistencies in segmentation outcomes due to ambiguous descriptions or variations in interpretation. Placing visual prompts directly on the image is often more efficient than composing detailed text prompts, especially for complex medical images like MRI or CT scans. This approach allows for better adaptation to the complexity and diversity of such images. However, the effectiveness of visual prompts depends on the user's ability to accurately place them, which requires a certain level of professional expertise. Additionally, visual prompts are not easily scalable for batch segmentation of medical images, as each target may have a different location. Text prompts, by contrast, are often integrated into expert models for specific tasks and offer unique value in particular scenarios. These prompts can take various forms, such as reports, templates, or labels, and usually require an encoder to convert them into vectors before input into the segmentation model. On one hand, text-prompt-based models can be trained with multiple types and modalities of data simultaneously, making the segmentation model more versatile. On the other hand, text-prompt-based models can be trained with various types and modalities of data simultaneously, making the segmentation model more versatile and capable of segmenting different organs and lesions. On the other hand, text prompts can capture anatomical relationships, improving segmentation performance and enhancing transfer learning on novel tasks.

Medical Image Classification

In medical imaging, classification tasks typically include disease classification, identifying benign and malignant tumors, molecular subtyping, lesion detection, and nucleus detection and classification. Prompt-based methods for medical image classification are primarily focused on foundational models, particularly CLIP, where classification is often performed as a downstream task. These foundational models are trained using contrastive learning, which has proven to be a highly effective and scalable strategy⁶⁸.

As a leading foundational model, CLIP¹⁷⁰ has indirectly driven advancements in medical image classification. Several medical-specific variants of CLIP^{21,22,62,182,183} have been developed, with an illustrative figure adapted from²², as shown in figure. 3. These models demonstrate robust zero-shot transfer capabilities through prompts, allowing for strong classification performance with only brief text prompts describing the target classes. As¹⁸⁴ confirmed, performance can be further enhanced by carefully designing appropriate text prompts.

Beyond foundational models, some expert models for medical image classification have also incorporated prompt mechanisms to improve performance. For example,⁶⁸ developed a prompt template and class name pool, randomly sampling 50 prompts to predict three image categories—introducing the first zero-shot transfer framework for histopathology whole-slide image classification. Similarly,⁶⁷ incorporated expert domain knowledge through descriptive textual prompts, enriching the limited categorical supervision typically found in medical datasets. Additionally,⁶⁵ used GPT-4 to acquire language-based prior knowledge at both the instance and bag levels, effectively addressing the challenge of Few-shot Weakly Supervised Whole-Slide Image Classification. In another approach,¹⁵⁴ used multiple prompts to describe medical lesions, enabling the detection of zero-shot lesions.

Moreover, several studies focus on learnable prompts.¹¹ achieved leading performance by introducing learnable visual prompt tokens. In contrast,⁶⁹ designed group prompts as learnable parameters to avoid the inefficiencies of fine-tuning the backbone for nucleus detection and classification.⁵ proposed a novel framework called Prompt-driven Latent Domain Generalization to address domain generalization in medical image classification without explicitly relying on domain labels, using a set of learnable token prompts. Similarly,⁷¹ introduced a Dynamic Prompt Module, which dynamically adapts intermediate representations to different dataset sources by leveraging dataset names and label properties as learnable prompts. This module uniformly predicts the corresponding categories of pathological images across various datasets.

Overall, prompt design is flexible and diverse, ranging from text and labels to learnable prompts. Using prompts helps alleviate the challenge of small data volumes while also improving model generalizability and adaptability. However, designing effective text prompts remains challenging, as prompt designs often lack interpretability. Additionally, foundational models may require substantial computational resources for training and fine-tuning.

Medical Image Foundation Models

As an essential method for interacting with foundation models, prompts play a crucial role in unlocking the potential of these models. This section summarizes the applications of medical foundation models in medical imaging in recent years. Foundation models—the latest generation of AI models—are trained on massive, diverse datasets and can be applied to numerous downstream tasks¹⁸⁵. Driven by growing datasets, increases in model size, and advances in model architectures, foundation models offer previously unseen abilities that promise to solve more diverse and challenging tasks than current

medical AI models, even while requiring little to no labels for specific tasks¹⁸⁶. To provide a comprehensive perspective, this section highlights the three most widely used applications of foundation models in medical imaging: histopathology, radiology, and ophthalmology.

Histopathology image evaluation is indispensable for cancer diagnoses and subtype classification¹⁸⁷. Pathology foundation models typically utilize contrastive learning methods, leveraging large numbers of whole-slide images for self-supervised learning. These pre-trained models can extract pathology imaging features for systematic cancer evaluation. Representative works in this area include references^{187, 188, 188, 189, 189, 190}. These foundation models have significantly outperformed state-of-the-art deep learning methods in tasks such as cancer detection, tumor origin characterization, genomic mutation identification, and survival prediction.

Compared to pathology images, radiology images present unique challenges for training foundation models due to the diverse formats and dimensions of data across different imaging modalities. Most foundation models in radiology are developed based on a single imaging modality. For example, in Chest X-ray imaging, representative foundation models such as MedCLIP²², CheXzero¹⁹¹, CXR-CLIP¹⁹², UniChest¹⁹³, and MedKLIP¹⁹⁴ have been developed for tasks such as disease classification. However, developing foundation models is highly challenging for 3D radiology data such as CT and MRI. Researchers often split volumes into slices or sub-volumes^{195, 196}. Nevertheless, some studies^{99, 182, 183, 197–199} have explored 3D foundation models for tasks such as zero-shot findings classification, phenotype classification, zero-shot cross-modal retrieval, disease prediction, radiology report generation, and segmentation. Although these models demonstrate significant progress, they also highlight that comprehensive 3D foundation models for radiology are still being developed.

Another widely used type of foundation model in medicine is the Ophthalmology foundation model. Not only are these models^{13, 200–202} used to diagnose common conditions such as diabetic retinopathy, glaucoma, and age-related macular degeneration, but they can also reveal indicators of systemic conditions like early-onset Parkinson's disease and cardiovascular disorders. In addition to fundus photography and OCT, other imaging modalities such as angiography, slit-lamp imaging, and ultrasound are also used to develop foundation models^{203, 204}. These models support a wider range of tasks and outperform expert-designed models.

Future Directions and Open Challenges

The advent of prompt-based mechanisms signifies a paradigm shift from traditional, often static, AI models towards more dynamic, interactive, and context-aware frameworks in medical imaging. As illustrated in figure 8, while traditional deep learning models often operate as isolated, task-specific entities, prompt-driven AI evolves towards systems that can flexibly adapt to the complex demands of diverse medical imaging tasks and dynamic clinical environments. This inherent dynamism and the ability to integrate contextual information enhance model performance, generalization, and interpretability. However, realizing the full potential of prompt-driven AI in medicine necessitates addressing several open challenges and exploring promising future research trajectories.

Enhancing Model Capabilities through Advanced Prompting

Future research will focus on evolving prompt mechanisms to further unlock the capabilities of underlying AI models, particularly large foundation models.

Towards Unified Multi-Task and Multi-Modal Generalists

The complexity of real-world medical practice demands AI systems capable of handling many tasks and data types¹⁸⁶. Prompt engineering is pivotal in transitioning from single-task, unimodal models to unified multi-task generalists. For example, a single foundation model, guided by different prompts, could perform target detection, image synthesis, and segmentation in various organs and imaging modalities (e.g., CT, MRI, PET), as demonstrated by early efforts such as seq2seq models for synthesis and segmentation³⁹ and UMS-Rep for shared knowledge across tasks²⁰⁵. Future work should aim to develop more sophisticated prompt-based architectures that enhance cross-task knowledge transfer and parameter efficiency, enabling robust performance across a broader spectrum of clinical scenarios.

Concurrently, the transition from unimodal to multimodal prompting holds immense potential⁵³. Clinical decision making inherently relies on integrating information from various sources (e.g., images, text-based electronic health records, genomic data, and patient history). Multimodal prompts, which condition a model on inputs from several modalities simultaneously, can provide richer contextual information than unimodal prompts, leading to superior performance. For example, TSF-Seq2Seq utilizes image and text information from other modalities as prompts to synthesize images of a specified target modality²⁰⁶. SimTxtSeg leverages the cross-modal fusion of images and text to generate high-quality pseudo-labels in segmentation tasks⁵¹. Future research should explore more effective strategies to combine and align multimodal prompt information, developing models that can seamlessly reason across and integrate these diverse data streams to improve diagnostic accuracy and treatment planning.

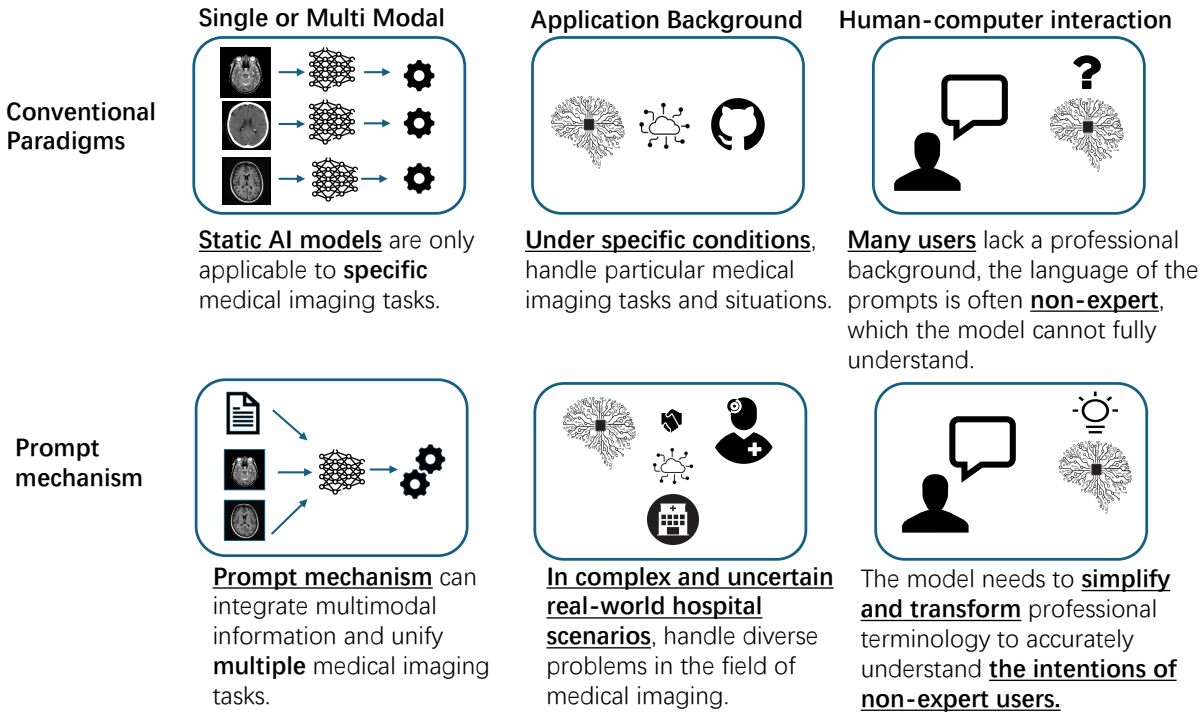


Figure 8. The future development of traditional deep learning and the future development of prompt-based deep learning.

Designing Superior, Robust, and Standardized Prompting Methodologies

The efficacy of current prompt-based systems heavily relies on the quality and design of the input prompts^{39,68,206}. However, existing prompt engineering practices face significant challenges:

- **High Sensitivity and Fragility:** Models often exhibit high sensitivity to minor variations in prompt phrasing or structure, often due to the vast, high-dimensional input space they operate in and the specific patterns they learned during pre-training, impacting output consistency and reliability.
- **Lack of Standardization and Automation:** Prompt design is frequently an empirical, manual, and experience-driven process, lacking unified optimization standards, robust automated generation techniques, or systematic evaluation frameworks, leading to variability in performance.
- **Limited Support for Complex Reasoning Tasks:** Current prompts often struggle to effectively guide models through complex, multi-stage reasoning or multi-role interactions, as seen in the need for specialized techniques like Chain-of-Thought prompting, potentially leading to information loss or misinterpretation.

Addressing these limitations is crucial. Inspired by advances like CoCoOp²⁰⁷, which enhances prompt generalization by learning conditional input tokens, and SegVol⁴⁶, which demonstrates segmentation of numerous anatomical categories using semantic and spatial prompts on a large scale, future research must focus on developing more robust, generalizable, and standardized prompting methods. This includes creating automated prompt optimization techniques, establishing benchmarks for prompt evaluation, developing formalisms for complex task decomposition via prompts, and designing prompts less susceptible to adversarial perturbations or unintended biases.

Seamless Clinical Integration and Practical Deployment

For prompt-driven AI to make a tangible impact, seamless integration into existing clinical workflows and practical deployment considerations are paramount.

Responsive and Adaptive Clinical Prompting Systems

Current medical ecosystems involve diverse information systems (e.g., EMR, PACS) storing heterogeneous data (including structured data such as lab results and vital signs, alongside unstructured clinical notes and reports). To effectively integrate

prompt models, it is essential to overcome system incompatibility and data diversification challenges to enable real-time analysis and interpretation that supports clinical decision making²⁰⁸. Future efforts should therefore focus on building highly responsive and adaptive prompt system frameworks. This involves developing architectures capable of processing complex, real-time clinical data streams, dynamically generating precise prompts tailored to evolving clinical scenarios, and ensuring seamless, interoperable connections with existing medical information systems to reduce data silos.

Hardware-Agnostic and Efficient Model Deployment

The scalability of prompt-based models for deployment on existing hospital hardware, which may be computationally constrained despite being expensive, is critical. Running large foundation models, even with efficient prompting, can be challenging. Therefore, it is necessary to develop effective model scaling and compression techniques-such as knowledge distillation^{209–211}, which transfers knowledge from a large model to a smaller one, weight quantization²¹², parameter sharing^{213,214}, and efficient attention mechanisms²¹⁵ tailored to prompted FMs. These methods should enable dynamic adaptation to diverse computational environments, facilitating efficient inference on resource-limited devices prevalent in clinical settings.

Optimizing Human–AI Collaboration in Clinical Workflows

Prompt models can substantially improve collaboration efficiency between clinicians and AI systems. Using human-provided prompts, AI can summarize complex medical data or highlight critical findings, thereby supporting clinical diagnosis. However, determining the optimal division of labor between humans and AI remains a significant challenge. Studies have shown that suboptimal human-AI interaction can sometimes result in performance that is worse than AI operating alone²¹⁶; nevertheless, AI assistance often provides significant benefits for junior clinicians²¹⁷. Future work should focus on designing intuitive and efficient prompting interfaces and interaction protocols, with optimizations for different clinical tasks (e.g., real-time surgical guidance versus image-based diagnostic review). This requires moving beyond generic interaction methods, such as simple clicks or bounding-box selections common in SAM, and toward developing multi-modal prompting tools (e.g., text, voice, images, masks) that optimize human–AI interaction and enhance performance in clinical settings.

Advancing Human-Computer Interaction for Broader Accessibility

Significant advances in human-computer interaction are required to make prompt-driven AI accessible and effective for a wider range of users, including those without deep technical or medical experience.

Natural Language Interaction for Non-Expert Users

Current prompt models primarily rely on LLMs to process specialized natural language. However, in real-world medical settings, users (such as patients or administrative staff) often lack professional backgrounds and use lay expressions. Consequently, models must be strengthened in their ability to understand, interpret, and respond to non-expert language. Future research should focus on optimizing conversational AI to convert complex medical terminology into accessible language, accurately capture the intent behind non-expert queries, and deliver practical guidance clearly and intuitively. This approach can draw on systems such as Mani-GPT²¹⁸ and DiagGPT²¹⁹, which emphasize advanced natural interaction and dynamic dialogue management, including context tracking and user intent recognition. The ultimate goal is to broaden the user base and enhance user experience across a wider range of medical applications through more intuitive, adaptive dialogue strategies while ensuring robust model performance.

Co-Adaptive Learning and Personalized Prompts

Future systems will evolve towards co-adaptive learning, creating a synergistic loop where the AI and its user learn from one another. On one hand, the system will personalize its outputs by understanding a user's interaction patterns, such as adapting to a radiologist who prefers concise, finding-focused summaries over one requiring detailed anatomical descriptions²²⁰. On the other hand, more effective prompts optimize the human-AI partnership. This powerful adaptive capability extends beyond user preferences to patient-specific data; by generating highly personalized prompts based on an individual's medical history, genetic predispositions, and ongoing treatments, these systems can deliver precise, tailored support for diagnosis and treatment planning, truly embodying the principles of personalized medicine²²¹.

Conclusion

Prompt-based mechanisms represent a transformative advance in medical imaging, significantly enhancing deep learning models' adaptability, precision, and utility, especially for large foundation models. By enabling the effective integration of contextual information and expert knowledge through textual, visual, or multimodal prompts, these methodologies are pivotal in addressing long-standing challenges such as limited annotated data, data heterogeneity, and the complexity of medical scenarios. This review has systematically surveyed the landscape, highlighting how prompts facilitate more accurate image generation, precise delineation of anatomical structures in segmentation, and robust and interpretable classification, even in challenging

zero-shot and few-shot learning settings. The continued synergy between advanced prompt engineering and the evolving capabilities of foundation models holds immense potential to revolutionize medical imaging. This will not only advance diagnostic accuracy and optimize treatment outcomes but will also enhance the scalability and accessibility of cutting-edge AI-powered healthcare systems globally, ultimately contributing to more efficient, equitable, and patient-centered medicine.

References

1. Mattjie, C. *et al.* Zero-shot performance of the segment anything model (sam) in 2d medical imaging: A comprehensive evaluation and practical guidelines. In *2023 IEEE 23rd International Conference on Bioinformatics and Bioengineering (BIBE)*, 108–112 (IEEE Computer Society, 2023).
2. Akrou, M. *et al.* Diffusion-based data augmentation for skin disease classification: Impact across original medical datasets to fully synthetic images. *arXiv preprint arXiv:2301.04802* (2023).
3. Pinaya, W. H. *et al.* Brain imaging generation with latent diffusion models. In *MICCAI Workshop on Deep Generative Models*, 117–126 (Springer, 2022).
4. Fischer, M., Bartler, A. & Yang, B. Prompt tuning for parameter-efficient medical image segmentation. *Med. Image Analysis* **91**, 103024 (2024).
5. Yan, S. *et al.* Prompt-driven latent domain generalization for medical image classification. *IEEE Transactions on Med. Imaging* (2024).
6. Chambon, P., Bluethgen, C., Langlotz, C. P. & Chaudhari, A. Adapting pretrained vision-language foundational models to medical imaging domains. *arXiv preprint arXiv:2210.04133* (2022).
7. Liu, J. *et al.* Clip-driven universal model for organ segmentation and tumor detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 21152–21164 (2023).
8. Chen, Z., Xu, Q., Liu, X. & Yuan, Y. Un-sam: Universal prompt-free segmentation for generalized nuclei images. *arXiv preprint arXiv:2402.16663* (2024).
9. Gao, Y., Xia, W., Hu, D. & Gao, X. Desam: Decoupling segment anything model for generalizable medical image segmentation. *arXiv preprint arXiv:2306.00499* (2023).
10. Cheng, D. *et al.* Sam on medical images: A comprehensive study on three prompt modes. *arXiv preprint arXiv:2305.00035* (2023).
11. Zhang, Y. *et al.* Text-guided foundation model adaptation for pathological image classification. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 272–282 (Springer, 2023).
12. Zhu, Y., Li, K., Yu, L. & Heng, P. A. Memory-efficient prompt tuning for incremental histopathology classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, 7802–7810 (2024).
13. Silva-Rodriguez, J., Chakor, H., Kobbi, R., Dolz, J. & Ayed, I. B. A foundation language-image model of the retina (flair): Encoding expert knowledge in text supervision. *Med. Image Analysis* **99**, 103357 (2025).
14. Ye, Y., Xie, Y., Zhang, J., Chen, Z. & Xia, Y. Uniseg: A prompt-driven universal segmentation model as well as a strong representation learner. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 508–518 (Springer, 2023).
15. Cao, Q., Xu, Z., Chen, Y., Ma, C. & Yang, X. Domain prompt learning with quaternion networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 26637–26646 (2024).
16. Zhou, N. *et al.* Medsam-u: Uncertainty-guided auto multi-prompt adaptation for reliable medsam. *arXiv preprint arXiv:2409.00924* (2024).
17. Singhal, K. *et al.* Large language models encode clinical knowledge. *arXiv preprint arXiv:2212.13138* (2022).
18. Wahd, A. S. *et al.* Sam2rad: A segmentation model for medical images with learnable prompts. *arXiv preprint arXiv:2409.06821* (2024).
19. Yellapragada, S. *et al.* Pathldm: Text conditioned latent diffusion model for histopathology. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 5182–5191 (2024).
20. Xu, Y. *et al.* Medsyn: Text-guided anatomy-aware synthesis of high-fidelity 3d ct images. *IEEE Transactions on Med. Imaging* (2024).
21. Zhang, S. *et al.* Large-scale domain-specific pretraining for biomedical vision-language processing. *arXiv preprint arXiv:2303.00915* (2023).

22. Wang, Z., Wu, Z., Agarwal, D. & Sun, J. Medclip: Contrastive learning from unpaired medical images and text. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 3876–3887 (2022).
23. Gu, Y. *et al.* Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Comput. for Healthc. (HEALTH)* **3**, 1–23 (2021).
24. Huang, K., Altosaar, J. & Ranganath, R. Clinicalbert: Modeling clinical notes and predicting hospital readmission. *arXiv preprint arXiv:1904.05342* (2019).
25. Huang, Z. *et al.* Ccnet: Criss-cross attention for semantic segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, 603–612 (2019).
26. Chen, C.-F. R., Fan, Q. & Panda, R. Crossvit: Cross-attention multi-scale vision transformer for image classification. In *Proceedings of the IEEE/CVF international conference on computer vision*, 357–366 (2021).
27. Petit, O. *et al.* U-net transformer: Self and cross attention for medical image segmentation. In *Machine Learning in Medical Imaging: 12th International Workshop, MLMI 2021, Held in Conjunction with MICCAI 2021, Strasbourg, France, September 27, 2021, Proceedings 12*, 267–276 (Springer, 2021).
28. Lin, H., Cheng, X., Wu, X. & Shen, D. Cat: Cross attention in vision transformer. In *2022 IEEE international conference on multimedia and expo (ICME)*, 1–6 (IEEE, 2022).
29. Dai, L., Zhang, R., Huang, Z. & Zhang, X. Guidegen: A text-guided framework for joint ct volume and anatomical structure generation. *arXiv preprint arXiv:2403.07247* (2024).
30. Bluethgen, C. *et al.* A vision–language foundation model for the generation of realistic chest x-ray images. *Nat. Biomed. Eng.* 1–13 (2024).
31. Chambon, P. *et al.* Roentgen: vision-language foundation model for chest x-ray generation. *arXiv preprint arXiv:2211.12737* (2022).
32. Hashmi, A. U. R. *et al.* Xreal: Realistic anatomy and pathology-aware x-ray generation via controllable diffusion model. *arXiv preprint arXiv:2403.09240* (2024).
33. Liang, Z., Xue, Z., Rajaraman, S. & Antani, S. Covid-19 pneumonia chest x-ray pattern synthesis by stable diffusion. In *2024 IEEE Southwest Symposium on Image Analysis and Interpretation (SSIAI)*, 21–24 (IEEE, 2024).
34. Han, W., Kim, C., Ju, D., Shim, Y. & Hwang, S. J. Advancing text-driven chest x-ray generation with policy-based reinforcement learning. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 56–66 (Springer, 2024).
35. Liu, C., Yuan, X., Yu, Z. & Wang, Y. Texdc: Text-driven disease-aware 4d cardiac cine mri images generation. In *Proceedings of the Asian Conference on Computer Vision*, 3005–3021 (2024).
36. Hofmanninger, J. *et al.* Automatic lung segmentation in routine imaging is primarily a data diversity problem, not a methodology problem. *Eur. Radiol. Exp.* **4**, 1–13 (2020).
37. Wang, A., Tam, T. C. C., Poon, H. M., Yu, K.-C. & Lee, W.-N. Naviairway: a bronchiole-sensitive deep learning-based airway segmentation pipeline. *arXiv preprint arXiv:2203.04294* (2022).
38. Wasserthal, J. *et al.* Totalsegmentator: robust segmentation of 104 anatomic structures in ct images. *Radiol. Artif. Intell.* **5** (2023).
39. Han, L. *et al.* Non-adversarial learning: Vector-quantized common latent space for multi-sequence mri. In Linguraru, M. G. *et al.* (eds.) *Medical Image Computing and Computer Assisted Intervention – MICCAI 2024*, 481–491 (Springer Nature Switzerland, Cham, 2024).
40. Han, L. *et al.* Synthesis-based imaging-differentiation representation learning for multi-sequence 3d/4d mri. *Med. Image Analysis* **92**, 103044 (2024).
41. Moghadam, P. A. *et al.* A morphology focused diffusion probabilistic model for synthesis of histopathology images. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2000–2009 (2023).
42. Ma, J. *et al.* Segment anything in medical images. *Nat. Commun.* **15**, 654 (2024).
43. Li, H., Liu, H., Hu, D., Wang, J. & Oguz, I. Promise: Prompt-driven 3d medical image segmentation using pretrained image foundation models. In *2024 IEEE International Symposium on Biomedical Imaging (ISBI)*, 1–5 (IEEE, 2024).
44. Deng, G. *et al.* Sam-u: Multi-box prompts triggered uncertainty estimation for reliable sam in medical image. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 368–377 (Springer, 2023).

45. Wang, G. *et al.* Sam-med3d-moe: Towards a non-forgetting segment anything model via mixture of experts for 3d medical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 552–561 (Springer, 2024).
46. Du, Y., Bai, F., Huang, T. & Zhao, B. Segvol: Universal and interactive volumetric medical image segmentation. *arXiv preprint arXiv:2311.13385* (2023).
47. Ramesh, D. B., Iytha Sridhar, R., Upadhyaya, P. & Kamaleswaran, R. Lugsam: A novel framework for integrating text prompts to segment anything model (sam) for segmentation tasks of icu chest x-rays. *Pulakesh Kamaleswaran, Rishikesan, Lugsam: A Nov. Framew. for Integrating Text Prompts to Segm. Anything Model. (Sam) for Segmentation Tasks Icu Chest X-Rays*.
48. Koleilat, T., Asgariandehkordi, H., Rivaz, H. & Xiao, Y. Medclip-samv2: Towards universal text-driven medical image segmentation. *arXiv preprint arXiv:2409.19483* (2024).
49. Kirillov, A. *et al.* Segment anything. *arXiv preprint arXiv:2304.02643* (2023).
50. Zhong, Y., Xu, M., Liang, K., Chen, K. & Wu, M. Ariadne’s thread: Using text prompts to improve segmentation of infected areas from chest x-ray images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 724–733 (Springer, 2023).
51. Xie, Y., Zhou, T., Zhou, Y. & Chen, G. Simtxtseg: Weakly-supervised medical image segmentation with simple text cues. *arXiv preprint arXiv:2406.19364* (2024).
52. Li, W., Xiong, X., Xia, P., Ju, L. & Ge, Z. Tp-drseg: Improving diabetic retinopathy lesion segmentation with explicit text-prompts assisted sam. *arXiv preprint arXiv:2406.15764* (2024).
53. Ye, Y., Chen, Z., Zhang, J., Xie, Y. & Xia, Y. Meduniseg: 2d and 3d medical image segmentation via a prompt-driven universal model. *arXiv preprint arXiv:2410.05905* (2024).
54. Liu, J. *et al.* Dctp-net: Dual-branch clip-enhance textual prompt-aware network for acute ischemic stroke lesion segmentation from ct image. *IEEE J. Biomed. Heal. Informatics* (2024).
55. Lin, L. *et al.* Fedlppa: Learning personalized prompt and aggregation for federated weakly-supervised medical image segmentation. *arXiv preprint arXiv:2402.17502* (2024).
56. Lin, Y. *et al.* Multi-target domain adaptation with prompt learning for medical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 717–727 (Springer, 2023).
57. Na, S., Guo, Y., Jiang, F., Ma, H. & Huang, J. Segment any cell: A sam-based auto-prompting fine-tuning framework for nuclei segmentation. *arXiv preprint arXiv:2401.13220* (2024).
58. Luo, W. *et al.* Universal medical image segmentation with task-specific prompt-guided transformer model. In *2023 International Annual Conference on Complex Systems and Intelligent Science (CSIS-IAC)*, 569–575 (IEEE, 2023).
59. Chen, Z., Pan, Y., Ye, Y., Lu, M. & Xia, Y. Each test image deserves a specific prompt: Continual test-time adaptation for 2d medical image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 11184–11193 (2024).
60. Zhang, S. *et al.* Biomedclip: a multimodal biomedical foundation model pretrained from fifteen million scientific image-text pairs. *arXiv preprint arXiv:2303.00915* (2023).
61. Eslami, S., Meinel, C. & De Melo, G. Pubmedclip: How much does clip benefit visual question answering in the medical domain? In *Findings of the Association for Computational Linguistics: EACL 2023*, 1181–1193 (2023).
62. Eslami, S., de Melo, G. & Meinel, C. Does clip benefit visual question answering in the medical domain as much as it does in the general domain? *arXiv preprint arXiv:2112.13906* (2021).
63. Bie, Y., Luo, L., Chen, Z. & Chen, H. Xcoop: Explainable prompt learning for computer-aided diagnosis via concept-guided context optimization. *arXiv preprint arXiv:2403.09410* (2024).
64. Han, M. *et al.* Mscpt: Few-shot whole slide image classification with multi-scale and context-focused prompt tuning. *arXiv preprint arXiv:2408.11505* (2024).
65. Qu, L., Fu, K., Wang, M., Song, Z. *et al.* The rise of ai language pathologists: Exploring two-level prompt learning for few-shot weakly-supervised whole slide image classification. *Adv. Neural Inf. Process. Syst.* **36** (2024).
66. Chikontwe, P., Kang, M., Luna, M., Nam, S. & Park, S. H. Low-shot prompt tuning for multiple instance learning based histology classification. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 285–295 (Springer, 2024).

67. Silva-Rodriguez, J., Chakor, H., Kobbi, R., Dolz, J. & Ayed, I. B. A foundation language-image model of the retina (flair): Encoding expert knowledge in text supervision. *arXiv preprint arXiv:2308.07898* (2023).
68. Lu, M. Y. *et al.* Visual language pretrained multiple instance zero-shot transfer for histopathology images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 19764–19775 (2023).
69. Huang, J., Li, H., Sun, W., Wan, X. & Li, G. Prompt-based grouping transformer for nucleus detection and classification. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 569–579 (Springer, 2023).
70. Zhu, W. *et al.* Segprompt: Using segmentation map as a better prompt to finetune deep models for kidney stone classification. In *Medical Imaging with Deep Learning*, 1680–1690 (PMLR, 2024).
71. Huang, J., Li, H., Wan, X. & Li, G. Unicell: Universal cell nucleus classification via prompt learning. *arXiv preprint arXiv:2402.12938* (2024).
72. Ye, Y., Zhang, J. & Shi, H. Pseudo-prompt generating in pre-trained vision-language models for multi-label medical image classification. In *Chinese Conference on Pattern Recognition and Computer Vision (PRCV)*, 279–298 (Springer, 2024).
73. Lin, Y., Zhu, Z., Cheng, K.-T. & Chen, H. Prompt-guided adaptive model transformation for whole slide image classification. *arXiv preprint arXiv:2403.12537* (2024).
74. Sánchez Quijada, M. Exploring large vision-language models with prompt engineering for peripheral blood cell image analysis and classification. *Univ. Oberta de Catalunya (UOC)* (2024).
75. Hamamci, I. E. *et al.* Generatect: Text-conditional generation of 3d chest ct volumes. In *European Conference on Computer Vision*, 126–143 (Springer, 2025).
76. Wang, Y. *et al.* Towards general text-guided image synthesis for customized multimodal brain mri generation. *arXiv preprint arXiv:2409.16818* (2024).
77. Shi, S., Li, H., Zhang, Y. & Wang, X. Semantic information-guided attentional gan-based ultrasound image synthesis method. *Biomed. Signal Process. Control.* **102**, 107273 (2025).
78. Dahan, E. *et al.* Csg: A context-semantic guided diffusion approach in de novo musculoskeletal ultrasound image generation. *arXiv preprint arXiv:2412.05833* (2024).
79. Yu, Y. *et al.* Ct synthesis with conditional diffusion models for abdominal lymph node segmentation. *arXiv preprint arXiv:2403.17770* (2024).
80. Xiao, Q. & Zhao, L. End-to-end 3d liver ct image synthesis from vasculature using a multi-task conditional generative adversarial network. *Appl. Sci.* **13**, 6784 (2023).
81. Weber, T., Ingrisch, M., Bischl, B. & Rügamer, D. Cascaded latent diffusion models for high-resolution chest x-ray synthesis. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, 180–191 (Springer, 2023).
82. Shentu, J. & Al Moubayed, N. Cxr-irgen: An integrated vision and language model for the generation of clinically accurate chest x-ray image-report pairs. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 5212–5221 (2024).
83. Borghesi, A. & Calegari, R. Generation of clinical skin images with pathology with scarce data. In *AI for Health Equity and Fairness: Leveraging AI to Address Social Determinants of Health*, 47–64 (Springer, 2024).
84. Fang, Z. *et al.* Conditional diffusion model for x-ray segmentation data generation. *J. Artif. Intell. Pract.* **7**, 7–10 (2024).
85. Sagers, L. W. *et al.* Improving dermatology classifiers across populations using images generated by large diffusion models. *arXiv preprint arXiv:2211.13352* (2022).
86. Wang, Y. *et al.* Toward general text-guided multimodal brain mri synthesis for diagnosis and medical image analysis. *Cell Reports Medicine* (2025).
87. Li, L. *et al.* Interactive gadolinium-free mri synthesis: A transformer with localization prompt learning. *arXiv preprint arXiv:2503.01265* (2025).
88. Duan, Y. *et al.* Fetalflex: Anatomy-guided diffusion model for flexible control on fetal ultrasound image synthesis. *arXiv preprint arXiv:2503.14906* (2025).
89. Liu, Q. *et al.* Treatment-aware diffusion probabilistic model for longitudinal mri generation and diffuse glioma growth prediction. *IEEE Transactions on Med. Imaging* (2025).

90. Wu, J. & Xu, M. One-prompt to segment all medical images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 11302–11312 (2024).
91. Chang, A. *et al.* Pe-med: Prompt enhancement for interactive medical image segmentation. In *International Workshop on Machine Learning in Medical Imaging*, 257–266 (Springer, 2023).
92. Bai, F. *et al.* Slpt: Selective labeling meets prompt tuning on label-limited lesion segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 14–24 (Springer, 2023).
93. Wu, C., Restrepo, D., Shuai, Z., Liu, Z. & Shen, L. Efficient in-context medical segmentation with meta-driven visual prompt selection. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 255–265 (Springer, 2024).
94. Xu, Y., Tang, J., Men, A. & Chen, Q. Eviprompt: A training-free evidential prompt generation method for segment anything model in medical images. *arXiv preprint arXiv:2311.06400* (2023).
95. Xie, B., Tang, H., Duan, B., Cai, D. & Yan, Y. Masksam: Towards auto-prompt sam with mask classification for medical image segmentation. *arXiv preprint arXiv:2403.14103* (2024).
96. Kato, S. & Hotta, K. One-shot and partially-supervised cell image segmentation using small visual prompt. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4295–4304 (2023).
97. Zhang, Y. *et al.* Continual learning for abdominal multi-organ and tumor segmentation. In *International conference on medical image computing and computer-assisted intervention*, 35–45 (Springer, 2023).
98. Tomar, N. K., Jha, D., Bagci, U. & Ali, S. Tganet: Text-guided attention for improved polyp segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 151–160 (Springer, 2022).
99. Zhao, Z. *et al.* One model to rule them all: Towards universal segmentation for medical images with text prompts. *arXiv preprint arXiv:2312.17183* (2023).
100. Biswas, R. Polyp-sam++: Can a text guided sam perform better for polyp segmentation? *arXiv preprint arXiv:2308.06623* (2023).
101. Chen, Y., Wang, Y. & Xie, Z. Vp-sfda: Visual prompt source-free domain adaptation for cross-modal medical image. *Heal. Data Sci.* .
102. Han, X., Chen, Q., Xie, Z., Li, X. & Yang, H. Multiscale progressive text prompt network for medical image segmentation. *Comput. & Graph.* **116**, 262–274 (2023).
103. Saeed, N., Ridzuan, M., Majzoub, R. A. & Yaqub, M. Prompt-based tuning of transformer models for multi-center medical image segmentation of head and neck cancer. *Bioengineering* **10**, 879 (2023).
104. Li, X., Zhang, Y. & Zhao, L. Multi-prompt fine-tuning of foundation models for enhanced medical image segmentation. *arXiv preprint arXiv:2310.02381* (2023).
105. Xie, B., Tang, H., Cai, D., Yan, Y. & Agam, G. Self-prompt sam: Medical image segmentation via automatic prompt sam adaptation. *arXiv preprint arXiv:2502.00630* (2025).
106. Xu, Q. *et al.* Sppnet: A single-point prompt network for nuclei image segmentation. In Cao, X., Xu, X., Rekik, I., Cui, Z. & Ouyang, X. (eds.) *Machine Learning in Medical Imaging*, 227–236 (Springer Nature Switzerland, Cham, 2024).
107. Sridhar, R. I. & Kamaleswaran, R. Lung segment anything model (lusam): A prompt-integrated framework for automated lung segmentation on icu chest x-ray images. *Authorea Prepr.* (2023).
108. Zhang, S., Yue, J., Wang, C., Liu, X. & Wang, G. Box2pseudo: A semi-supervised learning framework for pulmonary nodule segmentation with box-prompt pseudo supervision. In *2023 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, 1696–1703 (IEEE, 2023).
109. Glatt, R. & Shusen, L. Topology data analysis guided prompt optimization of segment anything model for zero-shot segmentation of biological images. Tech. Rep., Lawrence Livermore National Laboratory (LLNL), Livermore, CA (United States) (2023).
110. Zhou, Q., Feng, Y., Huang, Z., Ding, M. & Zhang, X. Specific instance and cross prompt based robust 3d semi-supervised medical image segmentation. *IEEE Transactions on Instrumentation Meas.* (2024).
111. Huang, Y. *et al.* Robust box prompt based sam for medical image segmentation. In *International Workshop on Machine Learning in Medical Imaging*, 1–11 (Springer, 2024).
112. Ouyang, X. *et al.* Prompt-based segmentation model of anatomical structures and lesions in ct images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 522–532 (Springer, 2024).

113. Chen, Y. *et al.* Segmentation by registration-enabled sam prompt engineering using five reference images. In *International Workshop on Biomedical Image Registration*, 241–252 (Springer, 2024).
114. Shaharabany, T. & Wolf, L. Zero-shot medical image segmentation based on sparse prompt using finetuned sam. In *Medical Imaging with Deep Learning* (2024).
115. Wang, R., Zhuang, L., Chen, H., Xu, B. & Cai, R. Tp-unet: Temporal prompt guided unet for medical image segmentation. *arXiv preprint arXiv:2411.11305* (2024).
116. Adhikari, R., Thapaliya, S., Dhakal, M. & Khanal, B. Tunevlseg: Prompt tuning benchmark for vision-language segmentation models. In *Proceedings of the Asian Conference on Computer Vision*, 126–144 (2024).
117. Kong, Y., Kim, K., Jeong, S., Lee, K. E. & Kong, H. Swiftmedsam: An ultra-lightweight prompt-based universal medical image segmentation model for highly constrained environments. In *CVPR 2024: Segment Anything In Medical Images On Laptop*.
118. Liu, X. *et al.* Feature-prompting gbmseg: One-shot reference guided training-free prompt engineering for glomerular basement membrane segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 276–285 (Springer, 2024).
119. Chen, H. *et al.* Multi-organ foundation model for universal ultrasound image segmentation with task prompt and anatomical prior. *IEEE Transactions on Med. Imaging* (2024).
120. Cui, C. *et al.* Enhancing physician flexibility: Prompt-guided multi-class pathological segmentation for diverse outcomes. In *IEEE-EMBS International Conference on Biomedical and Health Informatics*.
121. Xie, J. *et al.* Promamba: Prompt-mamba for polyp segmentation. *arXiv preprint arXiv:2403.13660* (2024).
122. Xia, Y. *et al.* Cervical-yosa: Utilizing prompt engineering and pre-trained large-scale models for automated segmentation of multi-sequence mri images in cervical cancer. *IET Image Process.* **18**, 3556–3569 (2024).
123. Teng, L. *et al.* Knowledge-guided prompt learning for lifespan brain mr image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 238–248 (Springer, 2024).
124. Chen, Z. *et al.* Adaptation of prompt-enabled segment-anything-model enhance the accuracy and generalizability of cine cardiac magnetic resonance segmentation. *Circulation* **150**, A4143921–A4143921 (2024).
125. Guan, H., Dai, B. & Zhang, J. Lite class-prompt tiny-vit for multi-modality medical image segmentation. In Ma, J., Zhou, Y. & Wang, B. (eds.) *Medical Image Segmentation Foundation Models. CVPR 2024 Challenge: Segment Anything in Medical Images on Laptop*, 151–166 (Springer Nature Switzerland, Cham, 2025).
126. Song, J., Yun, S., Yoon, S., Kim, J. & Lee, S. Ep-sam: Weakly supervised histopathology segmentation via enhanced prompt with segment anything. *arXiv preprint arXiv:2410.13621* (2024).
127. Khor, H. G. *et al.* Unified prompt-visual interactive segmentation of clinical target volume in ct for nasopharyngeal carcinoma with prior anatomical information. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 659–669 (Springer, 2024).
128. Xue, X. *et al.* Deep learning-based segmentation for high-dose-rate brachytherapy in cervical cancer using 3d prompt-resunet. *Phys. Medicine & Biol.* **69**, 195008 (2024).
129. Cui, C. *et al.* All-in-sam: from weak annotation to pixel-wise nuclei segmentation with prompt-based finetuning. In *Journal of Physics: Conference Series*, vol. 2722, 012012 (IOP Publishing, 2024).
130. Lyu, F., Xu, J., Zhu, Y., Wong, G. L.-H. & Yuen, P. C. Superpixel-guided segment anything model for liver tumor segmentation with couinaud segment prompt. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 678–688 (Springer, 2024).
131. Yang, J., Huang, Y., He, X., Shen, L. & Qiu, G. Tavp: Task-adaptive visual prompt for cross-domain few-shot segmentation. *arXiv preprint arXiv:2409.05393* (2024).
132. Xue, X. *et al.* A deep learning-based 3d prompt-nnunet model for automatic segmentation in brachytherapy of postoperative endometrial carcinoma. *J. Appl. Clin. Med. Phys.* e14371 (2024).
133. Dai, P., Ou, Y., Yang, Y., Liu, Y. & Zhao, Y. Sparse anatomical prompt semi-supervised learning with masked image modeling for cbct tooth segmentation. In *2024 IEEE International Symposium on Biomedical Imaging (ISBI)*, 1–5 (IEEE, 2024).
134. Cui, C. *et al.* Pfps: Prompt-guided flexible pathological segmentation for diverse potential outcomes using large vision and language models. *arXiv preprint arXiv:2407.09979* (2024).

135. Hu, K. & Xu, C. Lpm: A lightweight medical segmentation network based on mamba improved by prompt attention. *IET Image Process.* **18**, 3545–3555 (2024).
136. Sun, Y., Liu, M. & Lian, C. Aepl: Automated and editable prompt learning for brain tumor segmentation. *arXiv preprint arXiv:2410.19847* (2024).
137. Song, Y., Zhang, Y. & Li, M. An automatic laryngoscopic image segmentation system based on sam prompt engineering: From glottis annotation to vocal fold segmentation. *Authorea Prepr.* (2024).
138. Cheng, Y. & Zheng, Y. Frequency filtering prompt tuning for medical image semantic segmentation with missing modalities. *Big Data Inf. Anal.* **8**, 109–128 (2024).
139. Li, Y., Ren, H., Deng, J., Ma, X. & Xie, X. Centersam: Fully automatic prompt for dense nucleus segmentation. In *2024 IEEE International Symposium on Biomedical Imaging (ISBI)*, 1–5 (IEEE, 2024).
140. Zhang, Q., Guo, H., Yang, S., Li, Q. & Wang, Y. Progressive vision-language prompt for multi-organ multi-class cell semantic segmentation with single branch. *arXiv preprint arXiv:2412.02978* (2024).
141. Huang, X., He, D., Li, Z., Zhang, X. & Wang, X. Iossam: Label efficient multi-view prompt-driven tooth segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 632–642 (Springer, 2024).
142. Li, W. *et al.* Btsspro: Prompt-guided multimodal co-learning for breast cancer tumor segmentation and survival prediction. *IEEE J. Biomed. Heal. Informatics* (2024).
143. Shan, D. *et al.* Stpnet: Scale-aware text prompt network for medical image segmentation. *IEEE Transactions on Image Process.* (2025).
144. Wang, H. *et al.* Weakmedsam: Weakly-supervised medical image segmentation via sam with sub-class exploration and prompt affinity mining. *IEEE Transactions on Med. Imaging* (2025).
145. Yin, D., Zheng, Q., Chen, L., Hu, Y. & Wang, Q. Apg-sam: Automatic prompt generation for sam-based breast lesion segmentation with boundary-aware optimization. *Expert. Syst. with Appl.* **276**, 127048 (2025).
146. Liu, S., Zhang, D. & Hao, X. Efficient deformable convolutional prompt for continual test-time adaptation in medical image segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 39, 5550–5557 (2025).
147. Yin, S., Liu, S. & Wang, M. Ddfp: Data-dependent frequency prompt for source free domain adaptation of medical image segmentation. *Knowledge-Based Syst.* 113651 (2025).
148. Gao, Y. *et al.* Dual-prompt-enhanced multiorgan segmentation model for total-body pet images. *IEEE Transactions on Radiat. Plasma Med. Sci.* (2025).
149. Tian, C. *et al.* Self-prompt contextual learning with axialmamba for multi-label segmentation in carotid ultrasound. *Expert. Syst. with Appl.* **274**, 126749 (2025).
150. Zou, J. *et al.* Acea-net: Weakly supervised prostate 3d mri image segmentation via advanced prompt points. *IEEE J. Biomed. Heal. Informatics* (2025).
151. Chen, Z., Xu, Q., Liu, X. & Yuan, Y. Un-sam: Domain-adaptive self-prompt segmentation for universal nuclei images. *Med. Image Analysis* 103607 (2025).
152. Zhang, Y. *et al.* Category prompt mamba network for nuclei segmentation and classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 39, 10284–10292 (2025).
153. Zhao, J. *et al.* Uncertainty-driven edge prompt generation network for medical image segmentation. *IEEE Transactions on Med. Imaging* (2025).
154. Guo, M. *et al.* Multiple prompt fusion for zero-shot lesion detection using vision-language models. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 283–292 (Springer, 2023).
155. Cao, Q., Xu, Z., Chen, Y., Ma, C. & Yang, X. Domain prompt learning with quaternion networks. *arXiv preprint arXiv:2312.08878* (2023).
156. Zheng, F. *et al.* Exploring low-resource medical image classification with weakly supervised prompt learning. *Pattern Recognit.* **149**, 110250 (2024).
157. Huang, Y., Cheng, P., Tam, R. & Tang, X. Fine-grained prompt tuning: A parameter and memory efficient transfer learning method for high-resolution medical image classification. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 120–130 (Springer, 2024).

158. Yang, L. & Qu, W. Using text-augmented visual prompt learning for histopathology image classification. In *2024 5th International Conference on Big Data & Artificial Intelligence & Software Engineering (ICBASE)*, 272–276 (IEEE, 2024).
159. Bai, Y., Bai, L., Yang, X. & Liang, J. Label-semantic-based prompt tuning for vision transformer adaptation in medical image analysis. *IEEE Transactions on Circuits Syst. for Video Technol.* (2025).
160. Koleilat, T., Asgariandehkordi, H., Rivaz, H. & Xiao, Y. Biomedcoop: Learning to prompt for biomedical vision-language models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, 14766–14776 (2025).
161. He, A., Wu, Y., Wang, Z., Li, T. & Fu, H. Dvpt: Dynamic visual prompt tuning of large pre-trained models for medical image analysis. *Neural Networks* **185**, 107168 (2025).
162. Luo, Y. *et al.* Llm-guided decoupled probabilistic prompt for continual learning in medical image diagnosis. *IEEE Transactions on Med. Imaging* (2025).
163. Shin, H.-C. *et al.* Medical image synthesis for data augmentation and anonymization using generative adversarial networks. In *Simulation and Synthesis in Medical Imaging: Third International Workshop, SASHIMI 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16, 2018, Proceedings 3*, 1–11 (Springer, 2018).
164. Li, X. *et al.* Artificial general intelligence for medical imaging. *arXiv preprint arXiv:2306.05480* (2023).
165. Rombach, R., Blattmann, A., Lorenz, D., Esser, P. & Ommer, B. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10684–10695 (2022).
166. Goodfellow, I. *et al.* Generative adversarial networks. *Commun. ACM* **63**, 139–144 (2020).
167. Bowles, C. *et al.* Gan augmentation: Augmenting training data using generative adversarial networks. *arXiv preprint arXiv:1810.10863* (2018).
168. Kwon, G., Han, C. & Kim, D.-s. Generation of 3d brain mri using auto-encoding generative adversarial networks. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 118–126 (Springer, 2019).
169. Sun, L. *et al.* Hierarchical amortized gan for 3d high resolution medical image synthesis. *IEEE journal biomedical health informatics* **26**, 3966–3975 (2022).
170. Radford, A. *et al.* Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763 (PMLR, 2021).
171. Zhang, J., Xie, Y., Xia, Y. & Shen, C. Dodnet: Learning to segment multi-organ and tumors from multiple partially labeled datasets. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 1195–1204 (2021).
172. Zhao, Y. *et al.* Tg-lmm: Enhancing medical image segmentation accuracy through text-guided large multi-modal model. *arXiv preprint arXiv:2409.03412* (2024).
173. Boecking, B. *et al.* Making the most of text semantics to improve biomedical vision–language processing. In *European conference on computer vision*, 1–21 (Springer, 2022).
174. Mazurowski, M. A. *et al.* Segment anything model for medical image analysis: an experimental study. *Med. Image Analysis* **89**, 102918 (2023).
175. Hu, C. & Li, X. When sam meets medical images: An investigation of segment anything model (sam) on multi-phase liver tumor segmentation. *arXiv preprint arXiv:2304.08506* (2023).
176. Deng, R. *et al.* Segment anything model (sam) for digital pathology: Assess zero-shot segmentation on whole slide imaging. *arXiv preprint arXiv:2304.04155* (2023).
177. Roy, S. *et al.* Sam. md: Zero-shot medical image segmentation capabilities of the segment anything model. *arXiv preprint arXiv:2304.05396* (2023).
178. Cheng, J. *et al.* Sam-med2d. *arXiv preprint arXiv:2308.16184* (2023).
179. Wang, H. *et al.* Sam-med3d. *arXiv preprint arXiv:2310.15161* (2023).
180. Devlin, J. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
181. Chen, P., Li, Q., Biaz, S., Bui, T. & Nguyen, A. gscorecam: What objects is clip looking at? In *Proceedings of the Asian Conference on Computer Vision*, 1959–1975 (2022).
182. Hamamci, I. E. *et al.* A foundation model utilizing chest ct volumes and radiology reports for supervised-level zero-shot detection of abnormalities. *CoRR* (2024).

183. Lu, Z., Li, H., Parikh, N. A., Dillman, J. R. & He, L. Radclip: Enhancing radiologic image analysis through contrastive language-image pre-training. *arXiv preprint arXiv:2403.09948* (2024).
184. Zhou, K., Yang, J., Loy, C. C. & Liu, Z. Learning to prompt for vision-language models. *Int. J. Comput. Vis.* **130**, 2337–2348 (2022).
185. Bommasani, R. *et al.* On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258* (2021).
186. Moor, M. *et al.* Foundation models for generalist medical artificial intelligence. *Nature* **616**, 259–265 (2023).
187. Wang, X. *et al.* A pathology foundation model for cancer diagnosis and prognosis prediction. *Nature* **634**, 970–978 (2024).
188. Vorontsov, E. *et al.* A foundation model for clinical-grade computational pathology and rare cancers detection. *Nat. medicine* 1–12 (2024).
189. Xiang, J. *et al.* A vision–language foundation model for precision oncology. *Nature* 1–10 (2025).
190. Huang, Z., Bianchi, F., Yuksekogonul, M., Montine, T. J. & Zou, J. A visual–language foundation model for pathology image analysis using medical twitter. *Nat. medicine* **29**, 2307–2316 (2023).
191. Tiu, E. *et al.* Expert-level detection of pathologies from unannotated chest x-ray images via self-supervised learning. *Nat. Biomed. Eng.* **6**, 1399–1406 (2022).
192. You, K. *et al.* Cxr-clip: Toward large scale chest x-ray language-image pre-training. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 101–111 (Springer, 2023).
193. Dai, T. *et al.* Unichest: Conquer-and-divide pre-training for multi-source chest x-ray classification. *IEEE Transactions on Med. Imaging* (2024).
194. Wu, C., Zhang, X., Zhang, Y., Wang, Y. & Xie, W. Medklip: Medical knowledge enhanced language-image pre-training for x-ray diagnosis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 21372–21383 (2023).
195. Lei, Y., Li, Z., Shen, Y., Zhang, J. & Shan, H. Clip-lung: Textual knowledge-guided lung nodule malignancy prediction. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 403–412 (Springer, 2023).
196. Niu, C. *et al.* Medical multimodal-multitask foundation model for superior chest ct performance. *arXiv preprint arXiv:2304.02649* (2023).
197. Wu, C., Zhang, X., Zhang, Y., Wang, Y. & Xie, W. Towards generalist foundation model for radiology. *arXiv preprint arXiv:2308.02463* (2023).
198. Bai, F., Du, Y., Huang, T., Meng, M. Q.-H. & Zhao, B. M3d: Advancing 3d medical image analysis with multi-modal large language models. *arXiv preprint arXiv:2404.00578* (2024).
199. Blankemeier, L. *et al.* Merlin: A vision language foundation model for 3d computed tomography. *Res. Sq.* rs–3 (2024).
200. Zhou, Y. *et al.* A foundation model for generalizable disease detection from retinal images. *Nature* **622**, 156–163 (2023).
201. Engelmann, J. & Bernabeu, M. O. Training a high-performance retinal foundation model with half-the-data and 400 times less compute. *arXiv preprint arXiv:2405.00117* (2024).
202. Men, Y. *et al.* Drstagenet: Deep learning for diabetic retinopathy staging from fundus images. *arXiv preprint arXiv:2312.14891* (2023).
203. Qiu, J. *et al.* Visionfm: a multi-modal multi-task vision foundation model for generalist ophthalmic artificial intelligence. *arXiv preprint arXiv:2310.04992* (2023).
204. Shi, D. *et al.* Eyefound: A multimodal generalist foundation model for ophthalmic imaging. *arXiv preprint arXiv:2405.11338* (2024).
205. Zamzmi, G., Rajaraman, S. & Antani, S. Unified representation learning for efficient medical image analysis. *arXiv preprint arXiv:2006.11223* (2020).
206. Han, L. *et al.* An explainable deep framework: Towards task-specific fusion for multi-to-one mri synthesis. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 45–55 (Springer, 2023).
207. Zhou, K., Yang, J., Loy, C. C. & Liu, Z. Conditional prompt learning for vision-language models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 16816–16825 (2022).
208. Jiang, F. *et al.* Artificial intelligence in healthcare: past, present and future. *Stroke vascular neurology* **2** (2017).

209. Sinha, A. & Dolz, J. Multi-scale self-guided attention for medical image segmentation. *IEEE journal biomedical health informatics* **25**, 121–130 (2020).
210. Le Vuong, T. T. & Kwak, J. T. Moma: momentum contrastive learning with multi-head attention-based knowledge distillation for histopathology image analysis. *Med. Image Analysis* **101**, 103421 (2025).
211. Huang, K. *et al.* Learnable prompting sam-induced knowledge distillation for semi-supervised medical image segmentation. *IEEE Transactions on Med. Imaging* (2025).
212. Xie, L. *et al.* Mh-pflid: Model heterogeneous personalized federated learning via injection and distillation for medical data analysis. *arXiv preprint arXiv:2405.06822* (2024).
213. Shi, H., Ren, S., Zhang, T. & Pan, S. J. Deep multitask learning with progressive parameter sharing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 19924–19935 (2023).
214. Liu, P., Gao, Z.-F., Chen, Y., Zhao, W. X. & Wen, J.-R. Enhancing scalability of pre-trained language models via efficient parameter sharing. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, 13771–13785 (2023).
215. Qiu, Z. *et al.* Learning co-plane attention across mri sequences for diagnosing twelve types of knee abnormalities. *Nat. Commun.* **15**, 7637 (2024).
216. Kim, H.-E. *et al.* Changes in cancer detection and false-positive recall in mammography using artificial intelligence: a retrospective, multireader study. *The Lancet Digit. Heal.* **2**, e138–e148 (2020).
217. Tschandl, P. *et al.* Human–computer collaboration for skin cancer recognition. *Nat. medicine* **26**, 1229–1234 (2020).
218. Zhang, Z., Chai, W. & Wang, J. Mani-gpt: A generative model for interactive robotic manipulation. *Procedia Comput. Sci.* **226**, 149–156 (2023).
219. Cao, L. Diaggpt: An llm-based chatbot with automatic topic management for task-oriented dialogue. *arXiv preprint arXiv:2308.08043* (2023).
220. Shi, R. *et al.* From general to specific: Tailoring large language models for personalized healthcare. *arXiv preprint arXiv:2412.15957* (2024).
221. Shenfeld, I., Faltings, F., Agrawal, P. & Pacchiano, A. Language model personalization via reward factorization. *arXiv preprint arXiv:2503.06358* (2025).