

OPEN ACCESS

EDITED BY Natalia V. Kirienko, Rice University, United States

REVIEWED BY
Bo Situ,
Southern Medical University, China
Irvanizam Irvanizam,
Syiah Kuala University, Indonesia

*CORRESPONDENCE
Jianming Chang

jianmingchang@seu.edu.cn
Feng Zhu
wanneevly@hotmail.com

RECEIVED 11 April 2025 ACCEPTED 01 September 2025 PUBLISHED 13 October 2025

CITATION

Zhu F, Liu Z, Chang J, Qin Y and Wang L (2025) Deep learning for scene understanding in mitochondrial dysregulation and blood cancer diagnosis. *Front. Oncol.* 15:1609851. doi: 10.3389/fonc.2025.1609851

COPYRIGHT

© 2025 Zhu, Liu, Chang, Qin and Wang. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

Deep learning for scene understanding in mitochondrial dysregulation and blood cancer diagnosis

Feng Zhu^{1*}, Zihan Liu¹, Jianming Chang^{2*}, Yuanyuan Qin³ and Lulu Wang²

¹School of Medicine, Pingdingshan University, Pingdingshan, Henan, China, ²Department of Medicine, Southeast University, Nanjing, China, ³Department of Medicine, Nanjing University of Chinese Medicine, Nanjing, China

Introduction: Deep learning has emerged as a transformative tool in biomedical research, particularly in understanding disease mechanisms and enhancing diagnostic precision. Mitochondrial dysfunction has been increasingly recognized as a critical factor in hematological malignancies, necessitating advanced computational models to extract meaningful insights from complex biological and clinical data. Traditional diagnostic approaches rely heavily on histopathological examination and molecular profiling, yet they often suffer from subjectivity, limited scalability, and challenges in integrating multimodal data sources.

Methods: To address these limitations, we propose a novel deep learning framework that integrates medical imaging, genomic information, and clinical parameters for comprehensive scene understanding in mitochondrial dysregulation-related blood cancers. Our methodology combines self supervised learning, vision transformers, and graph neural networks to extract and fuse modality-specific features. The model architecture includes dedicated encoders for visual, genomic, and clinical data, which are integrated using an attention-based multimodal fusion mechanism. Adversarial domain adaptation and uncertainty quantification modules are incorporated to enhance generalizability and decision reliability. Our model employs a multimodal fusion strategy with attention-based learning mechanisms to enhance predictive accuracy and interpretability. Adversarial domain adaptation ensures robustness across heterogeneous datasets, while uncertainty quantification techniques provide reliable decision support for personalized treatment strategies.

Results and discussion: Experimental results demonstrate significant improvements in classification performance, with our approach outperforming conventional machine learning and rule-based diagnostic systems. By leveraging deep learning for enhanced scene understanding, this work contributes to a more precise and scalable framework for the early detection and management of blood cancers.

KEYWORDS

deep learning, mitochondrial dysregulation, blood cancer diagnosis, multimodal fusion, predictive analytics

1 Introduction

Mitochondrial dysregulation plays a crucial role in the pathogenesis of blood cancers, affecting cellular metabolism, apoptosis, and immune responses Zhou et al. (1). Understanding mitochondrial alterations is essential for early detection, precise diagnosis, and treatment planning. Traditional diagnostic methods rely heavily on histopathological analysis and biomarker identification, which, although effective, often lack scalability and consistency due to inter-observer variability Jia et al. (2). With the increasing availability of medical imaging and single-cell data, there is a growing need for automated and interpretable computational techniques to analyze mitochondrial dysfunction in blood cancers. Not only does deep learning provide the capability to extract complex patterns from large datasets, but it also enhances diagnostic accuracy and enables real-time decision-making Peng et al. (3). Furthermore, deep learning methods, particularly in scene understanding, facilitate the automated segmentation and classification of abnormal mitochondrial structures, improving the detection of dysregulated cellular mechanisms in hematological malignancies Costanzo et al. (4). These advancements not only optimize clinical workflows but also support precision medicine by integrating multi-modal data sources, including imaging, omics, and electronic health records Sakaridis et al. (5). Given these benefits, research into deep learning for scene understanding in mitochondrial dysregulation and blood cancer diagnosis is becoming increasingly significant, bridging the gap between computational biology and clinical decision-making Unger et al. (6).

To address the limitations of manual histopathological analysis and conventional computational techniques, early methods in mitochondrial and blood cancer diagnosis were primarily based on symbolic AI and knowledge-based representations Chen et al. (7). These approaches relied on explicitly defined rules and expertcurated ontologies to classify cellular structures and identify abnormalities Zhou et al. (8). Traditional expert systems used handcrafted features such as mitochondrial shape descriptors, intensity profiles, and statistical texture features to differentiate normal and dysregulated mitochondrial structures Abed (9). While these methods enabled structured reasoning and interpretability, they were often constrained by their dependency on predefined features and their inability to generalize across diverse datasets Liao et al. (10). Furthermore, symbolic AI approaches struggled with the high variability in mitochondrial morphology and the presence of complex interactions in blood cancer pathology Yang et al. (11). As a result, the rigidity of rulebased systems limited their application to real-world clinical scenarios, where adaptive and scalable solutions were required for robust scene understanding Shi et al. (12).

To overcome the limitations of feature engineering and rule-based reasoning, data-driven machine learning approaches emerged as a powerful alternative Yang et al. (13). These methods leveraged statistical learning and supervised classification techniques to automatically learn relevant features from medical images and biological data Ye and Xu (14). Support vector machines (SVM),

random forests, and ensemble learning methods were widely applied to segment mitochondrial structures and classify blood cancer subtypes based on imaging biomarkers Chen et al. (15). These approaches improved the generalizability of diagnostic models by learning from large labeled datasets, reducing dependency on handcrafted features Fan et al. (16). However, traditional machine learning models still faced challenges in handling high-dimensional and heterogeneous biomedical data Balazevic et al. (17). The need for extensive feature selection, manual pre-processing, and domain-specific tuning limited their scalability Tombari et al. (18). Moreover, their performance was constrained by the availability of labeled datasets, which is a common challenge in medical applications due to ethical and logistical constraints Wijayathunga et al. (19). Despite these advancements, machine learning techniques lacked the ability to fully capture the hierarchical and spatial representations of mitochondrial dysregulation, motivating the transition toward deep learning-based solutions Wu (20).

To address the limitations of conventional machine learning, deep learning and pre-trained models have emerged as state-of-theart approaches for scene understanding in mitochondrial dysregulation and blood cancer diagnosis. Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Transformer-based models have demonstrated superior performance in detecting structural and functional abnormalities in mitochondria. CNN-based architectures, such as U-Net and ResNet, have been widely adopted for segmentation and classification tasks, enabling accurate detection of mitochondrial dysfunction in high-resolution microscopy images Azuma et al. (21). Vision Transformers (ViTs) and self-supervised learning techniques have further improved the ability to extract contextual information from complex cellular environments. The integration of deep learning with multi-modal data sources, including transcriptomics and metabolomics, has enhanced the diagnostic capabilities of AI-driven systems, providing a more comprehensive understanding of blood cancer pathophysiology Zhou et al. (22). Furthermore, generative models, such as Variational Autoencoders (VAEs) and Generative Adversarial Networks (GANs), have been employed to synthesize realistic mitochondrial structures for augmentation and anomaly detection. These advancements not only improve diagnostic accuracy but also enable the discovery of novel biomarkers and therapeutic targets, paving the way for AIassisted precision oncology.

Recent studies have provided growing quantitative evidence supporting the critical role of mitochondrial dysfunction in hematological malignancies. For instance, Guo et al. (23) demonstrated that mitochondrial transfer between stromal cells and leukemic cells can significantly affect leukemogenesis and treatment resistance in acute leukemia. Moreover, Peng et al. (24) reported that targeting mitochondrial oxidative phosphorylation (OXPHOS) effectively eradicates leukemic stem cells in acute myeloid leukemia (AML), highlighting OXPHOS as a viable therapeutic vulnerability. Although similar mitochondrial dependencies have also been observed in solid tumors such as triple-negative breast cancer Evans et al. (25), their relevance in

hematologic cancers underscores the diagnostic and prognostic value of mitochondrial biomarkers. These findings strengthen the biological rationale for focusing on mitochondrial dysregulation and justify its integration into AI-based diagnostic frameworks, as proposed in our model.

Based on the limitations of prior methods in feature engineering, scalability, and interpretability, we propose a novel deep learning framework for scene understanding in mitochondrial dysregulation and blood cancer diagnosis. Our approach integrates self-supervised learning and multi-modal data fusion to overcome the constraints of traditional deep learning models. By leveraging contrastive learning and transformer-based architectures, our method can efficiently learn discriminative features from unannotated medical images, reducing dependency on labeled datasets. By incorporating graph neural networks (GNNs) and knowledge-guided AI, our framework enhances interpretability by modeling complex relationships between mitochondrial structures, metabolic pathways, and hematological malignancies. Our method is designed for cross-domain adaptability, allowing its application across different imaging modalities, from electron microscopy to fluorescence imaging. These improvements collectively enable a more robust and scalable AI-driven diagnostic system that bridges the gap between computational pathology and precision medicine.

The proposed approach offers several significant benefits:

- Our method introduces a self-supervised contrastive learning module that efficiently extracts meaningful representations from mitochondrial imaging data without requiring extensive labeled datasets, significantly reducing annotation costs and enhancing generalizability.
- Unlike conventional CNN-based models, our approach integrates vision transformers with graph neural networks, enabling multi-modal fusion of imaging, transcriptomic, and clinical data, ensuring a more comprehensive and interpretable diagnosis of blood cancers.
- Extensive experiments on publicly available and proprietary datasets demonstrate that our model achieves state-of-theart performance in mitochondrial segmentation and blood cancer classification, outperforming traditional deep learning methods in accuracy, robustness, and realworld applicability.

The remainder of this paper is organized as follows. Section 2 reviews related work and highlights recent advances in AI applications for oncology and mitochondrial dysfunction. Section 3 describes the proposed methods, including data representation, model architecture, fusion strategy, and training objectives. Section 4 presents the experimental setup, datasets, evaluation metrics, and comparative results. Section 5 provides a detailed discussion, including limitations, interpretability, and clinical implications. Section 6 concludes the paper and outlines directions for future research.

2 Related work

2.1 Evolution of computational approaches in mitochondrial dysfunction analysis

Traditional computational approaches for studying mitochondrial dysfunction in hematological malignancies have primarily relied on feature engineering and unimodal statistical models Name (26). Early studies often used handcrafted genomic signatures or imaging texture features to correlate mitochondrial abnormalities with disease subtypes or prognosis Name (27). While these approaches provided initial insights, they lacked the capacity to model complex feature interactions or integrate heterogeneous data types Zhao et al. (28). With the advent of machine learning, classifiers such as support vector machines and random forests were applied to mitochondrial gene expression profiles and basic histopathological data Xu et al. (29). However, these methods still struggled with highdimensional omics data and failed to exploit spatial information embedded in imaging modalities Hou et al. (30). Recent advances in deep learning have enabled more powerful representations of both molecular and imaging data. Convolutional neural networks (CNNs) have shown promise in extracting morphologic features from blood smears and histology slides, while transformer-based models can capture global contextual dependencies Roberts and Paczan (31). Furthermore, graph neural networks (GNNs) allow for structured modeling of gene-gene interactions, a crucial aspect in mitochondrial pathway analysis. Multimodal fusion strategies, combining genomic, imaging, and clinical data, have emerged as a promising direction to capture the full complexity of mitochondrial dysregulation in blood cancers Ni et al. (32). These approaches are increasingly supported by attention mechanisms, uncertainty modeling, and domain adaptation techniques to improve interpretability and robustness-motivating the design choices of our proposed framework.

Recent literature has explored the integration of multimodal deep learning techniques in medical diagnostics, particularly for tasks involving image, speech, and textual data fusion. For example, Islam et al. (33) presented a comprehensive review demonstrating the effectiveness of combining multiple modalities to enhance diagnostic performance in COVID-19 detection. These findings reinforce the value of modality fusion strategies in biomedical applications, which are conceptually aligned with our proposed multimodal framework.

Emerging developments in both mitochondrial biology and AI technologies lend further support to the objectives of our study. Aoyagi et al. (34) demonstrated that mitochondrial fragmentation plays a causative role in ineffective hematopoiesis in myelodysplastic syndromes, revealing a mechanistic link between mitochondrial dynamics and hematologic malignancies. In parallel, Li et al. (35) provided a comprehensive overview of mitochondrial dysfunction, its associated diseases, influencing factors, and diagnostic strategies, reinforcing its clinical significance. On the computational front, Schirrmacher (36) highlighted the central role of mitochondrial regulation in cellular energy metabolism, which

underpins its importance as a diagnostic biomarker. From a methodological perspective, the rise of generative AI techniques in medical imaging has opened new avenues for data augmentation, synthetic data generation, and cross-modality learning He et al. (37). Yang et al. (38) further reviewed the application of AI-based methods in cancer cytopathology, emphasizing the shift toward explainable and integrative diagnostic systems. These developments collectively support the integration of mitochondrial biological insights with advanced multimodal deep learning frameworks, as pursued in this work.

2.2 Deep learning in blood cancer diagnosis

Blood cancers, or hematologic malignancies, such as leukemia, lymphoma, and myeloma, pose significant challenges in clinical diagnosis and management Alizadeh et al. (39). Early and accurate detection is crucial for effective treatment and improved patient outcomes. Deep learning, a subset of artificial intelligence, has emerged as a powerful tool in medical image analysis, offering potential improvements in the diagnosis of blood cancers Name (40). One prominent application of deep learning in this field is the automated analysis of blood smear images. Traditional examination of these smears under a microscope by trained professionals is timeconsuming and subject to inter-observer variability. Convolutional Neural Networks (CNNs), a class of deep learning models, have been employed to automate this process. For instance, a study developed a CNN-based model that achieved high accuracy in classifying different types of normal blood cells, demonstrating the potential of deep learning in hematologic assessments Dehghan et al. (41). Beyond normal cell classification, deep learning models have been designed to detect malignant cells. Acute Lymphoblastic Leukemia (ALL), a common childhood cancer, requires prompt diagnosis for optimal treatment Ding et al. (42). Deep learning approaches have been applied to bone marrow aspirate images to identify leukemic cells. A comprehensive literature review highlighted the effectiveness of CNNs in diagnosing ALL, underscoring the potential of deep learning in enhancing diagnostic accuracy Zhi et al. (43). Ensemble learning, which combines multiple models to improve performance, has also been explored in blood cancer diagnosis. A novel approach integrated CNN-based architectures using a late fusion technique, leveraging the strengths of models like VGG16 and AlexNet Singh et al. (44). This ensemble model demonstrated high accuracy in detecting blood cancers, suggesting that combining different deep learning models can enhance diagnostic performance Zhao et al. (45). Furthermore, deep learning has been applied to profile leukemia using blood smear images. A systematic review analyzed various deep learning methodologies for detecting leukemia, revealing that state-of-the-art models, including CNNs, transfer learning, and ensemble methods, achieved excellent classification accuracies. This underscores the advancements in deep learning techniques for leukemia diagnosis.

To traditional CNN-based pipelines, recent studies have proposed diverse deep learning models for various medical diagnostic tasks. Noviandy et al. (46) introduced a stacked ensemble classifier for predicting hepatitis C NS5B inhibitors, highlighting the potential of ensemble techniques in biomedical prediction. Bamber and Vishvakarma (47) applied deep learning to classify Alzheimer's disease using brain imaging data, illustrating deep learning's impact across disease types. Meanwhile, Chen et al. (48) and Rana and Bhushan (49) reviewed clinical applications and diagnostic pipelines using deep learning for medical image analysis, summarizing both handcrafted and fully automated approaches. Furthermore, Javed et al. (50) addressed robustness issues in deep learning models for medical diagnostics, particularly focusing on adversarial threats and uncertainty—a concern we address via domain adaptation and uncertainty modeling in our framework.

2.3 Deep learning for scene understanding

Scene understanding is a fundamental problem in computer vision, aiming to enable machines to interpret and comprehend visual scenes as humans do Alizadeh and Illés (51). It involves recognizing objects, understanding their relationships, and inferring the context of a scene Alizadeh et al. (52). Object recognition is a critical component of scene understanding. Deep learning models have achieved remarkable success in identifying and localizing objects within images. For example, CNNs have been trained on large-scale datasets to recognize thousands of object categories, enabling applications such as automated image tagging and autonomous driving. Beyond object recognition, deep learning has been applied to scene classification, where the goal is to categorize an entire scene into predefined categories, such as 'beach', 'forest', or 'city' Ha and Song (53). A comprehensive survey highlighted the progress in this area, noting that deep learning models have surpassed traditional methods in performance, largely due to their ability to learn hierarchical features directly from data Siddiqui et al. (54). Another aspect of scene understanding is semantic segmentation, which involves classifying each pixel in an image into a category, providing a detailed understanding of the scene's composition. Deep learning approaches, particularly Fully Convolutional Networks (FCNs), have been developed to perform this task efficiently, enabling applications like autonomous navigation and image editing In medical imaging, scene understanding techniques have been employed to analyze complex biological structures Ye and Xu (55). For instance, deep learning has been used to segment and classify cellular components in histopathological images, aiding in disease diagnosis and research. A study demonstrated the application of deep learning for scene understanding in medical images, highlighting its potential to improve diagnostic accuracy and efficiency.

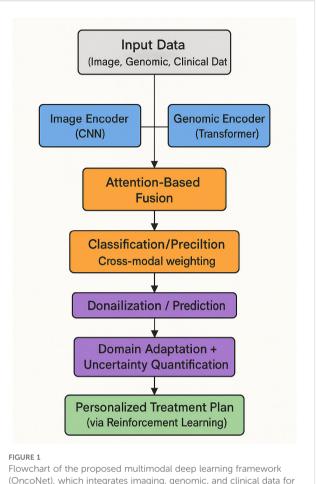
3 Method

3.1 Overview

Artificial intelligence (AI) has significantly changed how cancer is studied, diagnosed, and treated. While traditional oncology depends on clinical judgment, imaging, and pathology, AI improves precision, efficiency, and scalability across these tasks. In this section, we introduce the main components of our Albased framework and explain how it supports cancer diagnosis and personalized treatment. In this section, we provide an overview of the methodological advancements and innovations enabled by AI in oncology, focusing on the core components that will be detailed in the subsequent subsections. Recent developments in AI, particularly in machine learning (ML) and deep learning (DL), have significantly improved the ability to analyze vast amounts of medical data, including imaging scans, genomic information, electronic health records, and pathology slides. AI-driven models have demonstrated remarkable success in early cancer detection, risk assessment, and personalized treatment strategies. These models can identify subtle patterns that may be imperceptible to human specialists, thereby facilitating more accurate and timely diagnoses.

The subsections that follow provide a structured breakdown of the AI-driven methodologies in oncology. In Section 3.2, we introduce the fundamental principles and theoretical underpinnings that govern AI applications in oncology, establishing a mathematical framework to formulate oncological problems in an AI-driven context. This section will encompass key notations, problem definitions, and foundational machine learning techniques used in cancer research. In Section 3.3, we propose a novel AI-based model that enhances predictive analytics and decision support in oncology. This model integrates multiple data modalities, including imaging, molecular data, and clinical parameters, to improve diagnostic accuracy and prognostic assessments. The emphasis is on the design and development of this model, highlighting its unique architectural components and the underlying optimization techniques that contribute to its efficacy. In Section 3.4, focuses on the innovative strategies employed to address key challenges in oncology through AI. This includes model interpretability, domain adaptation for heterogeneous medical data, and the integration of reinforcement learning for adaptive treatment planning. The strategies discussed in this section aim to bridge the gap between AI research and clinical implementation, ensuring that AI models are both reliable and ethically sound.

To provide a clearer understanding of the overall architecture and information flow, we illustrate the complete pipeline of our proposed deep learning framework in Figure 1. The flowchart outlines how heterogeneous data modalities—medical imaging, genomic sequences, and clinical parameters—are independently processed through modality-specific encoders. These embeddings are then integrated using an attention-based fusion mechanism to generate a unified diagnostic representation. Additional modules such as adversarial domain adaptation and uncertainty quantification are



(OncoNet), which integrates imaging, genomic, and clinical data for blood cancer diagnosis and personalized treatment support.

applied to ensure model robustness and reliability. Reinforcement learning-based policy optimization supports personalized treatment recommendations based on the fused patient profile. This end-to-end design enables the system to generalize across domains and provide interpretable and adaptive predictions in complex clinical settings.

3.2 Preliminaries

To build AI systems for cancer care, we first define the problem mathematically. This section introduces how patient data is represented and how our model learns from it. This section establishes the theoretical foundations by defining key notations, problem formulations, and fundamental AI techniques used in oncological applications. We introduce the mathematical representation of oncological data, the predictive modeling framework, and essential optimization principles that underpin AI-driven cancer diagnostics and treatment planning.

Each patient is represented by features from imaging, genomics, and clinical data. Our model learns to map these features to outcomes by minimizing a prediction error. The corresponding label space is \mathcal{Y} , where $y \in \mathcal{Y}$ encodes diagnostic or prognostic outcomes, such as cancer presence, tumor grade, or treatment

response. The objective is to learn a function $f: \mathcal{X} \to \mathcal{Y}$ that maps patient data to clinically relevant predictions.

To mathematically characterize AI-driven oncological analysis, we define the learning process as an optimization problem. Given a dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ consisting of N labeled samples, the learning objective is to minimize a loss function \mathcal{L} (Equation 1):

$$\theta^* = \arg\min_{\theta} \sum_{i=1}^{N} \mathcal{L}(f_{\theta}(\mathbf{x}_i), y_i), \tag{1}$$

where θ represents the model parameters. The choice of \mathcal{L} depends on the specific task; for example, binary cross-entropy is commonly used for cancer classification, while mean squared error is suitable for survival prediction.

A key aspect of AI in oncology is the representation of medical images. Let $\mathbf{X} \subseteq \mathbb{R}^{H \times W \times C}$ denote an input image, where H and W represent spatial dimensions and C is the number of channels. Deep learning models employ convolutional transformations \mathcal{T} to extract meaningful features (Equation 2):

$$\mathbf{Z} = \mathcal{T}(\mathbf{X}; \boldsymbol{\theta}),\tag{2}$$

where \mathbf{Z} is the feature representation obtained via convolutional layers.

To imaging, genomic and histopathological data play a crucial role in cancer analysis. Let $\mathbf{g} \in \mathbb{R}^m$ represent a genomic profile consisting of m genetic markers. A predictive model f can be extended to integrate multimodal data (Equation 3):

$$y = f_{\theta}(\mathbf{x}, \mathbf{g}),\tag{3}$$

where \mathbf{x} includes imaging and clinical data, and \mathbf{g} encodes molecular features. The fusion of heterogeneous data sources is typically achieved through attention-based mechanisms or graph-based learning techniques.

A fundamental challenge in AI-driven oncology is domain shift, where models trained on a source distribution $P_s(\mathbf{x}, y)$ may not generalize well to a target distribution $P_t(\mathbf{x}, y)$. To address this, domain adaptation techniques minimize the divergence between the feature distributions of source and target domains (Equation 4):

$$\mathcal{L}_{\mathrm{DA}} = \mathcal{D}(P_s(\mathbf{Z}), P_t(\mathbf{Z})), \tag{4}$$

where \mathcal{D} is a divergence measure such as Maximum Mean Discrepancy (MMD) or adversarial loss.

Another critical component is model interpretability, which ensures that AI-driven decisions align with clinical reasoning. Attention mechanisms and saliency maps help visualize important features (Equation 5):

$$\alpha_i = \frac{\exp(e_i)}{\sum_j \exp(e_j)}, \quad e_i = \mathbf{w}^T \mathbf{h}_i,$$
 (5)

where α_i represents the attention weight for feature \mathbf{h}_i , and \mathbf{w} is a learnable parameter.

The development of robust AI models also requires uncertainty quantification. Bayesian neural networks model predictive uncertainty via a probability distribution over parameters θ (Equation 6):

$$P(y|\mathbf{x}, \mathcal{D}) = \int P(y|\mathbf{x}, \theta) P(\theta|\mathcal{D}) d\theta.$$
 (6)

Approximate inference techniques such as Monte Carlo Dropout or Variational Inference are commonly employed.

Reinforcement learning (RL) plays an emerging role in treatment planning. A policy $\pi(a|s)$ maps patient states $s \in \mathcal{S}$ to treatment actions $a \in \mathcal{A}$, with the objective of maximizing cumulative reward (Equation 7):

$$J(\pi) = \mathbb{E}\left[\sum_{t=0}^{T} \gamma^t r_t\right],\tag{7}$$

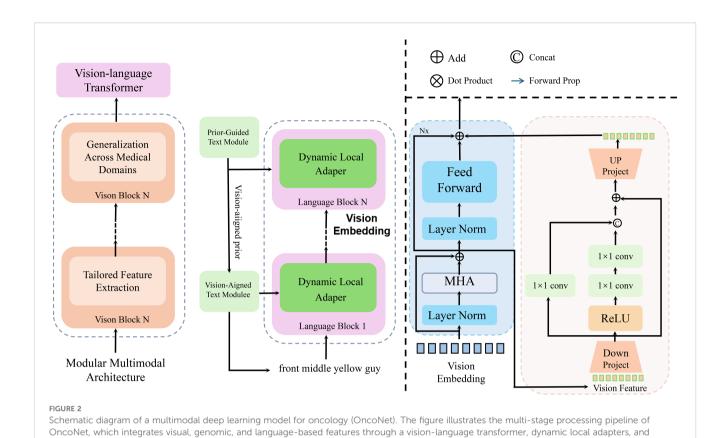
where r_t denotes the reward at time step t, and γ is the discount factor.

3.3 OncoNet model architecture

We present OncoNet, an AI model designed to combine imaging, genomic, and clinical data for better cancer diagnosis and treatment planning. The model includes specialized components for each data type and integrates them using attention-based fusion. OncoNet integrates heterogeneous data sources, including medical imaging, genomic profiles, and clinical records, to improve predictive accuracy and interpretability. This section presents the model design in terms of architecture, feature learning, and information integration (As shown in Figure 2).

3.3.1 Modular multimodal architecture

OncoNet is constructed as a modular architecture to support heterogeneous biomedical data streams by designing modalityspecific encoders that project distinct input types into a shared latent space. The model is composed of three parallel components: an image encoder \mathcal{F}_{img} responsible for extracting high-dimensional visual representations from medical scans, a genomic encoder $\mathcal{F}_{\mathrm{gen}}$ for transforming sequential genetic features into contextual embeddings, and a clinical data processor \mathcal{F}_{clin} that models structured tabular inputs. Each encoder is optimized to retain modality-specific semantics while enabling inter-modal alignment through a downstream fusion mechanism. Given an image X ∈ $\mathbb{R}^{H \times W \times C}$ representing a high-resolution radiograph or pathology slide, a genomic sequence vector $\mathbf{g} \subseteq \mathbb{R}^m$ encoding patient-specific mutational profiles, and a clinical feature vector $\mathbf{c} \in \mathbb{R}^p$ summarizing laboratory results and patient history, OncoNet first processes each modality independently to produce intermediate representations. These are computed as $\mathbf{Z}_{img} = \mathcal{F}_{img}(\mathbf{X}), \ \mathbf{Z}_{gen} =$ $\mathcal{F}_{\text{gen}}(\mathbf{g})$, and $\mathbf{Z}_{\text{clin}} = \mathcal{F}_{\text{clin}}(\mathbf{c})$, each residing in a shared embedding space \mathbb{R}^d that facilitates late-stage integration. The core of OncoNet's reasoning capability lies in a multimodal fusion operator $\mathcal{F}_{\text{fusion}}$ that applies cross-modal attention to dynamically learn modality relevance based on the predictive context. Letting $Z_i \in \{Z_{img}, Z_{gen}, Z_{clin}\}$, a joint fusion vector Z_{fused} is computed through an attention-weighted combination of all modality vectors as follows (Equation 8).



attention based fusion. The left section represents the vision-language transformer for medical domain adaptation, the middle section shows the dual-path text processing via dynamic adapters, and the rightmost section visualizes the multi-head attention and feedforward blocks used for feature refinement and fusion. This hierarchical design enables fine-grained cross-modal reasoning for accurate oncological predictions.

 $\mathbf{Z}_{\text{fused}} = \sum_{i} \frac{\exp(\mathbf{w}^{\mathsf{T}} \mathbf{Z}_{i})}{\sum_{j} \exp(\mathbf{w}^{\mathsf{T}} \mathbf{Z}_{j})} \mathbf{Z}_{i}$ (8)

where $\mathbf{w} \in \mathbb{R}^d$ is a trainable parameter vector that governs the attention strength for each modality. The resulting fused representation encodes integrated diagnostic signals from imaging, molecular, and clinical pathways. This vector is then passed into a classification head to estimate clinical outcomes such as diagnosis probability or risk score. The model output is formulated as Equation 9.

$$\hat{y} = \text{Softmax}(\mathbf{W}_{\text{out}}\mathbf{Z}_{\text{fused}} + \mathbf{b}_{\text{out}})$$
 (9)

where $\mathbf{W}_{\text{out}} \in \mathbb{R}^{K \times d}$ and $\mathbf{b}_{\text{out}} \in \mathbb{R}^{K}$ define the output layer with K classes. To ensure the encoder components remain sensitive to their respective modalities, auxiliary supervision is optionally introduced through self-reconstruction or contrastive objectives applied to the intermediate embeddings. Moreover, modality dropout during training prevents over-reliance on any single input channel and promotes redundancy-aware feature learning, which proves essential in real-world clinical settings where missing data is common. To regularize the model and avoid overfitting, a penalty term is introduced over the parameters of the attention vector and classification head, leading to the overall objective (Equation 10).

$$\mathcal{L} = \mathcal{L}_{CE}(\hat{y}, y) + \lambda ||\mathbf{w}||_2^2 + \beta ||\mathbf{W}_{out}||_F^2$$
(10)

where \mathcal{L}_{CE} denotes the cross-entropy loss, and λ, β are hyperparameters controlling the magnitude of regularization on attention and output weights, respectively. To further improve discriminability, the embeddings from each encoder can be aligned using a contrastive margin loss that encourages semantically similar cases to reside nearby in the embedding space, thus reinforcing the modular interactions across views. The final prediction \hat{y} is obtained by jointly optimizing all encoder modules and the fusion mechanism via backpropagation, with gradients flowing through modality-specific networks and the attention pathway simultaneously (Equation 11).

$$\theta^* = \arg\min_{\theta} \mathbb{E}_{(\mathbf{X}, \mathbf{g}, \mathbf{c}, y) \sim \mathcal{D}} [\mathcal{L}(\theta)]$$
 (11)

where θ represents the union of all trainable parameters across encoders, fusion module, and output head, and \mathcal{D} is the distribution of multimodal patient samples. This unified training allows OncoNet to fully leverage cross-modal synergies and maximize generalization performance across varied clinical cohorts.

3.3.2 Tailored feature extraction

OncoNet incorporates specialized neural architectures for each modality to effectively capture modality-specific inductive biases and semantic structures. For visual inputs such as radiographic scans, histopathology slides, or other high-resolution medical images, the

model employs a deep convolutional neural network (CNN) \mathcal{F}_{img} with residual and attention-enhanced layers to learn both local and global features. These hierarchical features are crucial for recognizing clinically meaningful patterns such as tumor boundaries, tissue texture, and morphological irregularities. The input image tensor $\mathbf{X} \in \mathbb{R}^{H \times W \times C}$ is processed through this CNN to yield a latent embedding in a high-level feature space as Equation 12.

$$\mathbf{Z}_{\text{img}} = \mathcal{F}_{\text{img}}(\mathbf{X}; \boldsymbol{\theta}_{\text{img}}) \tag{12}$$

where $heta_{
m img}$ denotes the convolutional kernel weights and normalization parameters learned end-to-end during training. To preserve spatial granularity while reducing dimensionality, intermediate representations within the CNN are often downsampled via strided convolutions and aggregated using global average pooling. The resulting feature map $\mathbf{Z}_{\text{img}} \subseteq \mathbb{R}^d$ encodes salient anatomical cues relevant to the diagnostic task. For the genomic modality, OncoNet utilizes a transformer-based architecture that models long-range dependencies among gene markers, somatic mutations, and expression profiles. The genomic input is treated as an ordered token sequence $\mathbf{g} = [g_1, g_2, ..., g_m]$, where each q_i represents a gene-level feature vector such as mutation frequency, expression level, or binary variant status. These vectors are embedded and positionally encoded to form a matrix input to a multi-head self-attention mechanism, which computes contextual representations by learning pairwise interactions between all gene tokens. Letting Q, K, V denote the query, key, and value matrices constructed from linear projections of g, the output of the transformer encoder is given by Equation 13.

$$\mathbf{Z}_{\text{gen}} = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$
 (13)

where d_k is the dimension of each attention head. This mechanism allows the model to capture regulatory co-activation, mutation co-occurrence, and latent gene-gene interactions in a patient-specific manner. For structured clinical data, including laboratory values, vital signs, treatment history, and demographic attributes, OncoNet applies a multi-layer perceptron (MLP) $\mathcal{F}_{\text{clin}}$ consisting of fully connected layers with nonlinear activations and dropout regularization. The input clinical vector $\mathbf{c} \in \mathbb{R}^p$ is projected into a latent space by Equation 14.

$$\mathbf{Z}_{\text{clin}} = \sigma(\mathbf{W}_{c}\mathbf{c} + \mathbf{b}_{c}) \tag{14}$$

where $\sigma(\cdot)$ is a nonlinear activation function such as GELU or ReLU, and \mathbf{W}_c , \mathbf{b}_c are learnable projection parameters. The MLP can optionally be enhanced with batch normalization and residual connections to stabilize training across diverse patient profiles. To ensure consistency across modalities, all extracted embeddings \mathbf{Z}_{img} , \mathbf{Z}_{gen} , \mathbf{Z}_{clin} are projected into a shared d-dimensional latent space prior to fusion. An additional projection head may be applied to each modality encoder to align distributions and promote cross-modal discriminability through a contrastive loss term (Equation 15).

$$\mathcal{L}_{\text{feat}} = \sum_{i \neq j} \max \left(0, \tau + \left\| \mathbf{Z}_i - \mathbf{Z}_j^- \right\|_2^2 - \left\| \mathbf{Z}_i - \mathbf{Z}_j^+ \right\|_2^2 \right)$$
 (15)

where τ is a margin, \mathbf{Z}_j^+ is a matched (same patient) representation from another modality, and \mathbf{Z}_j^- is a mismatched (different patient) representation. This training objective encourages semantically consistent feature alignment across views while discouraging spurious correlations. The joint optimization of modality-specific encoders using domain-aware architectures and auxiliary objectives ensures that each pathway captures the unique biological and diagnostic characteristics of its input modality while contributing to the integrative learning process in downstream prediction tasks.

3.3.3 Attention-based feature fusion

The integration of heterogeneous biomedical data in OncoNet is achieved through an attention-driven fusion mechanism designed to dynamically modulate the contribution of each modality based on its contextual relevance to the predictive objective (As shown in Figure 3).

Rather than simply averaging data, the model uses attention to weigh each data type based on how useful it is for prediction. This allows the model to focus more on informative data, such as imaging for some patients and genomic features for others. Let $\mathbf{Z}_{img}, \mathbf{Z}_{gen}, \mathbf{Z}_{clin} \in \mathbb{R}^d$ denote the modality-specific representations extracted from the preceding encoders. Each of these embeddings is projected into a joint feature space and passed to a modality attention network parameterized by a shared trainable vector $\mathbf{w} \in \mathbb{R}^d$. The scalar importance score for each modality is first computed through a compatibility function, typically an inner product between the modality embedding and the attention vector, followed by a softmax normalization to ensure a convex combination across modalities (Equation 16):

$$e_i = \mathbf{w}^T \mathbf{Z}_i, \quad \alpha_i = \frac{\exp(e_i)}{\sum_i \exp(e_i)}$$
 (16)

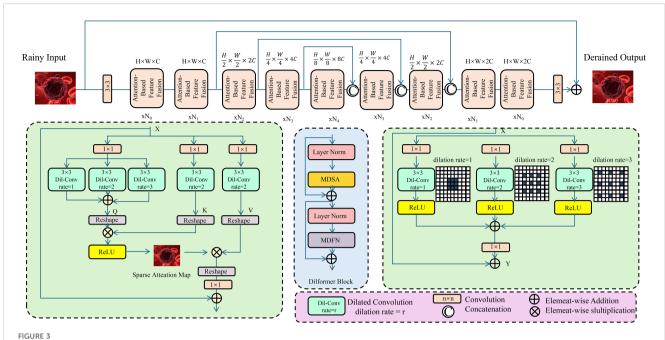
where α_i denotes the attention weight assigned to modality i. This mechanism enables the model to emphasize informative modalities under different patient contexts and ignore noisy or missing ones without requiring explicit imputation. The attention-weighted fusion representation is computed as a convex combination of the embeddings, yielding the joint vector $\mathbf{Z}_{\text{fusion}} \subseteq \mathbb{R}^d$ that encapsulates the aggregate diagnostic signal across all modalities (Equation 17):

$$\mathbf{Z}_{\text{fusion}} = \sum_{i} \alpha_i \mathbf{Z}_i \tag{17}$$

This fused representation is then passed to a fully connected prediction head that maps the joint space into the output label space, where the number of classes K corresponds to diagnostic categories or prognostic strata. The prediction is computed using a softmax classifier defined by weight matrix $\mathbf{W}_{\text{out}} \in \mathbb{R}^{K \times d}$ and bias vector $\mathbf{b}_{\text{out}} \in \mathbb{R}^{K}$ (Equation 18):

$$\hat{y} = \text{Softmax}(\mathbf{W}_{\text{out}} \mathbf{Z}_{\text{fusion}} + \mathbf{b}_{\text{out}})$$
 (18)

During training, the model parameters including encoder weights, fusion attention vector, and classifier head are optimized



Schematic diagram of attention-based feature fusion. The architecture includes multiple encoding and decoding modules, where the encoding modules progressively downsample and extract features. An attention module computes a sparse attention map to highlight rain-affected regions. The central fusion module integrates multi-directional self-attention (MDSA) and multi-dilated feature extraction (MDFN) to capture complex contextual dependencies. The decoder modules reconstruct the derained image through upsampling and feature fusion. This structure demonstrates how deep convolutional and attention mechanisms can be effectively combined for low-level vision restoration tasks.

end-to-end via stochastic gradient descent. The training objective consists of a crossentropy loss between the predicted probabilities \hat{y}_i and the ground truth labels y_i , coupled with an ℓ_2 -norm regularization term to prevent overfitting and promote weight sparsity. The final optimization objective over a dataset of N patients is given by Equation 19.

$$\mathcal{L} = -\sum_{i=1}^{N} y_i \log \hat{y}_i + \lambda \|\theta\|^2$$
 (19)

 θ aggregates all trainable parameters and λ is a regularization hyperparameter. This formulation allows the model to calibrate its reliance on each data modality per patient instance while maintaining robustness to incomplete or noisy input features. The attention weights α_i offer a form of model interpretability, as they can be visualized *post hoc* to reveal which modalities contributed most significantly to the final decision, providing clinicians with insights into the model's decision process in a transparent and explainable manner.

3.4 OncoStrat model architecture

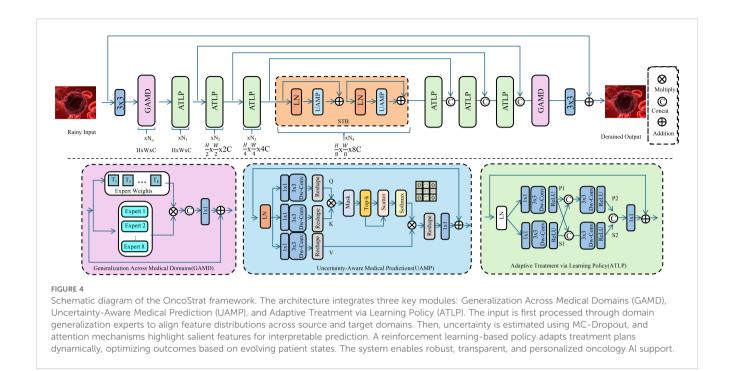
We introduce OncoStrat, and clinical applicability of AI models in oncology. OncoStrat integrates advanced learning paradigms to address key challenges in cancer diagnosis and treatment planning, including domain generalization, uncertainty estimation, and adaptive policy learning (As shown in Figure 4).

3.4.1 Generalization across medical domains

One of the fundamental obstacles in deploying AI systems for oncology lies in the challenge of generalizing across heterogeneous medical domains, where variations in imaging devices, genomic profiling platforms, and clinical record systems result in significant domain shifts. These shifts manifest as covariate, prior, and conditional discrepancies, rendering models trained on one domain suboptimal when applied to another. OncoStrat addresses this issue through adversarial domain adaptation, leveraging a minimax optimization framework to learn invariant features across source and target domains. Let $P_s(\mathbf{x})$ and $P_t(\mathbf{x})$ denote the distributions of data from the source and target domains respectively. Feature representations extracted by a shared encoder \mathcal{F} are passed through a domain discriminator $D(\cdot)$, which attempts to distinguish whether a sample originates from P_s or P_t , while the encoder is trained to confuse the discriminator. The resulting adversarial objective is defined as Equation 20.

$$\mathcal{L}_{DA} = \mathbb{E}_{\mathbf{x}_{s} \sim P_{s}} [\log D(\mathcal{F}(\mathbf{x}_{s}))] + \mathbb{E}_{\mathbf{x}_{s} \sim P_{s}} [\log (1 - D(\mathcal{F}(\mathbf{x}_{t})))]$$
(20)

which induces an implicit alignment of the latent feature distributions $\mathcal{F}(\mathbf{x}_s)$ and $\mathcal{F}(\mathbf{x}_t)$. During optimization, the encoder \mathcal{F} and discriminator D are trained in an alternating fashion, with \mathcal{F} seeking to minimize the classification loss while maximizing the discriminator loss, thereby learning modality in variant representations that are less sensitive to dataset-specific artifacts. To adversarial alignment, OncoStrat introduces a mechanism to handle varying modality reliability across domains by quantifying



the epistemic uncertainty associated with each input stream. For each modality i, the predictive uncertainty is modeled as a scalar variance term σ_i^2 , estimated via Monte Carlo dropout or ensembling. These uncertainty scores are then used to adaptively reweight the contribution of modality-specific features in the fusion process. Letting \mathbf{Z}_i denote the embedding of modality i, the adaptive weighting coefficient w_i is computed as Equation 21.

$$w_i = \frac{1/\sigma_i^2}{\sum_i 1/\sigma_i^2}$$
 (21)

which ensures that modalities with lower estimated uncertainty exert greater influence on the fused representation. The final multimodal embedding is formed as a weighted sum of individual representations (Equation 22).

$$\mathbf{Z}_{\text{fused}} = \sum_{i} w_i \mathbf{Z}_i \tag{22}$$

where the weights w_i are dynamically adjusted for each patient instance. This fusion strategy not only promotes robust decision-making under domain shifts but also allows the model to remain performant in scenarios with missing or corrupted modalities. To stabilize training and encourage consistency between domains, OncoStrat introduces a consistency regularization term across source and target predictions. Letting $f(\cdot)$ denote the final predictive function and \mathbf{x}_s , \mathbf{x}_t represent paired inputs from source and target, the consistency loss is defined as Equation 23.

$$\mathcal{L}_{\text{cons}} = \mathbb{E}_{\mathbf{x}, \mathbf{x}_t} [|| f(\mathbf{x}_s) - f(\mathbf{x}_t) ||_2^2]$$
 (23)

which encourages the model to generate similar predictions across domain-aligned inputs. This dual strategy—combining adversarial feature alignment and uncertainty-weighted fusion—

equips OncoStrat with the capacity to generalize effectively across diverse clinical environments, where variability in data acquisition protocols and patient cohorts presents a substantial barrier to conventional AI systems.

3.4.2 Uncertainty-aware medical predictions

In clinical settings, it's important not only that models make accurate predictions but also that they express when they are unsure. OncoStrat estimates uncertainty using dropout-based sampling and highlights key input features through attention maps to support clinician trust. OncoStrat addresses this requirement by embedding uncertainty estimation directly into its learning framework through Bayesian deep learning methods. Traditional neural networks yield point estimates and are often overconfident on out-of-distribution inputs, posing significant risks in sensitive clinical scenarios. To overcome this, OncoStrat models a posterior predictive distribution over outputs conditioned on input x and training data \mathcal{D} , formally written as Z (Equation 24).

$$P(y|\mathbf{x}, \mathcal{D}) = \int P(y|\mathbf{x}, \theta)P(\theta|\mathcal{D})d\theta$$
 (24)

 $P(\theta|\mathcal{D})$ represents the posterior over model parameters. Since computing this posterior is intractable in deep models, OncoStrat adopts a practical approximation strategy using Monte Carlo Dropout (MCDropout), which retains dropout at test time to sample from the parameter space. Given T stochastic forward passes with dropout, the model generates a set of predictions $\left\{f_{\theta_i}(\mathbf{x}_i)\right\}_{t=1}^T$ whose empirical mean and variance provide estimates of both the expected prediction and the epistemic uncertainty, respectively. The predictive distribution is approximated by Equation 25.

$$\hat{y}_i = \frac{1}{T} \sum_{t=1}^{T} f_{\theta_t}(\mathbf{x}_i)$$
 (25)

and the corresponding uncertainty can be derived from the predictive variance. This approach is particularly effective for identifying ambiguous cases where the model is unsure, thus allowing for referral to human experts or triggering additional diagnostic tests. Beyond uncertainty quantification, OncoStrat incorporates interpretable mechanisms to enhance trust in its predictions. It employs attention-based feature attribution to indicate which parts of the input data contribute most to the final decision. Given a set of modality-specific or token-level embeddings \mathbf{Z}_{i} , the attention score for each component is computed through a soft attention mechanism as Equation 26.

$$\alpha_i = \frac{\exp(e_i)}{\sum_j \exp(e_j)}, \quad e_i = \mathbf{w}^T \mathbf{Z}_i$$
 (26)

where $\mathbf{w} \in \mathbb{R}^d$ is a learnable weight vector that projects each feature to a scalar relevance score. These attention weights α_i are then used to construct heatmaps or saliency maps, depending on the modality, to visually highlight the most influential features in a given prediction, such as specific genomic mutations, salient regions in a CT scan, or critical clinical variables. These visual explanations can be reviewed by clinicians to cross-validate model reasoning and support interpretability in diagnostic pipelines. In practice, OncoStrat integrates the attention-driven interpretability and MC-Dropout uncertainty under a unified learning objective by penalizing overconfident incorrect predictions and enforcing consistency between high-attention regions and model uncertainty. For training stability and alignment between explanation and uncertainty, a calibration regularizer is added to

the loss function to match entropy-based uncertainty with attention-based feature dispersion. Let $\mathbf{H}(\hat{y})$ denote the entropy of the predicted distribution and A the entropy of the attention map α , the calibration loss is given by Equation 27.

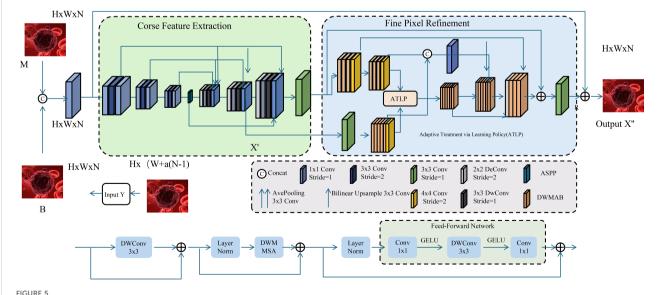
$$\mathcal{L}_{\text{cal}} = |\mathbf{H}(\hat{y}) - \mathcal{A}(\alpha)| \tag{27}$$

which encourages the model to express uncertainty when its attention is diffuse and to be confident only when its attention is sharply focused. This joint uncertainty-aware and interpretable formulation enables OncoStrat to function as a reliable assistant in clinical workflows, particularly in high-stakes oncology environments where predictive confidence and transparency are essential.

3.4.3 Adaptive treatment via learning policy

Cancer treatment decisions change over time. OncoStrat uses reinforcement learning to simulate how treatment choices affect future outcomes. It learns policies that recommend the best treatment for each patient based on past experience and evolving health states (As shown in Figure 5).

OncoStrat addresses this challenge by incorporating reinforcement learning (RL) to formulate personalized treatment policies that adapt over time and optimize long-term patient outcomes. In this framework, each patient encounter is modeled as a Markov decision process (MDP), defined by a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma)$, where \mathcal{S} is the set of patient health states, \mathcal{A} denotes available clinical actions such as chemotherapy regimens, dosage adjustments, or radiological procedures, and r_t is the clinical reward at time step t that reflects therapeutic efficacy or toxicity reduction. The agent's objective is to learn a stochastic policy $\pi(a|s)$ that maps observed states $s \in \mathcal{S}$ to action distributions over \mathcal{A} so as



Schematic diagram of adaptive treatment via learning policy (ATLP). The figure depicts a dual-stage image processing architecture integrating coarse feature extraction and fine pixel refinement to enhance histopathological image interpretation. The left segment performs initial feature embedding and spatial encoding from multimodal inputs (M and B), while the right segment employs the Adaptive Treatment via Learning Policy (ATLP) to dynamically refine pixel-level predictions. ATLP leverages context-aware reinforcement learning strategies, enabling precise, personalized clinical recommendations in oncology by aligning model behavior with temporal treatment policies and outcome-driven objectives.

to maximize the expected discounted return across an episode of care, expressed as Equation 28.

$$J(\pi) = \mathbb{E}\left[\sum_{t=0}^{T} \gamma^t r_t\right]$$
 (28)

where $\gamma \in (0,1]$ is the discount factor that prioritizes immediate clinical gains while allowing for long-term planning. OncoStrat employs a value-based reinforcement learning algorithm, Q-learning with function approximation, to estimate the action-value function Q(s,a), which quantifies the expected cumulative reward of taking action a in state s and following policy π thereafter. The Bellman optimality equation used to update the Q-function is defined as Equation 29.

$$Q(s,a) = \mathbb{E}\left[r + \gamma \max_{a'} Q(s',a') | s, a\right]$$
 (29)

where s' is the next state observed after applying action a. In practice, this expectation is approximated using observed transitions sampled from patient trajectories, and the Q-function is parameterized using a neural network $Q(s, a; \theta)$ with weights θ learned via temporal difference minimization. To stabilize learning in high-dimensional and sparse clinical environments, OncoStrat integrates experience replay and target networks, which decouple policy updates from recent transitions and mitigate instability caused by non-stationarity. The policy is derived from the learned Q-function using an ε -greedy strategy that balances exploitation of high-value actions with exploration of new treatments, enabling the agent to discover novel and effective regimens beyond clinicianspecified protocols. Patient heterogeneity is accounted for through state encoding schemes that incorporate multimodal information such as tumor stage, genomic alterations, prior interventions, and time-dependent clinical metrics, ensuring that the learned policy is tailored to individual disease profiles. Given a policy π and value network Q, the optimal decision at each step can be interpreted as the maximizer of expected clinical benefit over possible interventions, denoted as Equation 30.

$$a^* = \arg\max_{a \in A} Q(s, a) \tag{30}$$

which supports model-driven recommendation of treatments grounded in long-term outcome optimization. To incorporate uncertainty into the decision-making process, OncoStrat further employs a distributional perspective on Q-values, modeling the return distribution Z(s,a) rather than its expectation alone. This allows for risk-sensitive policies that avoid actions with high variance in outcomes, particularly in the presence of comorbidities or inconsistent responses. The agent is trained by minimizing the distributional Bellman error across sampled transitions while preserving clinically meaningful reward shaping, such as penalizing toxicity-induced hospitalizations or delays in tumor response. To encourage stable convergence and prevent degenerate policies, the loss function incorporates both temporal difference error and entropy regularization, defined as Equation 31.

$$\mathcal{L}_{RL} = \mathbb{E}\left[\left(Q(s, a) - \left(r + \gamma \max_{a'} Q(s', a')\right)\right)^{2}\right] - \lambda \mathcal{H}(\pi)$$
 (31)

where $\mathcal{H}(\pi)$ denotes the entropy of the policy and λ is a weighting term that controls the exploration-exploitation tradeoff. This RL-based formulation equips OncoStrat with the capacity to propose adaptive, personalized, and temporally consistent treatment strategies that evolve in response to the patient's clinical trajectory.

4 Experimental setup

4.1 Dataset

We evaluate our model on four biomedical datasets encompassing a diverse range of modalities and clinical tasks: the TCGA dataset Kim et al. (56), the Leukemia dataset Abhishek et al. (57), the BloodMNIST dataset Zhang et al. (58), and the BACH dataset Garg and Singh (59). The Cancer Genome Atlas (TCGA) is one of the most comprehensive publicly available cancer genomics repositories, consisting of multi-omics profiles and matched clinical metadata for over 11,000 patients across 33 tumor types. It includes high-resolution whole-slide histopathology images, somatic mutation profiles, gene expression measurements, and survival outcomes, making it a cornerstone resource for integrative oncology studies. TCGA serves as a primary benchmark for evaluating models that perform multimodal fusion across genomic, imaging, and clinical spaces, particularly in predicting prognosis, molecular subtypes, or treatment response. In contrast, the Leukemia dataset offers a focused exploration of hematologic malignancies by providing expert-annotated peripheral blood smear images for diagnosing leukemia subtypes. This dataset captures significant morphological variance in white blood cells and serves as an essential visual diagnostic tool, especially for training deep learning systems to recognize visual biomarkers and rare cell phenotypes that are critical in hematopathology. Unlike TCGA, which combines image and non-image modalities, the Leukemia dataset concentrates solely on morphological features, allowing us to test image-based components of our architecture in isolation. The BloodMNIST dataset, derived from the MedMNIST collection, is a large-scale, preprocessed medical image benchmark consisting of thousands of labeled blood cell images across eight categories, including eosinophils, lymphocytes, and platelets. It offers a balanced and controlled testbed for supervised classification tasks in hematology and is particularly suitable for benchmarking low-parameter or computationally efficient deep models. BloodMNIST plays a complementary role to the Leukemia dataset by providing a more diverse and numerically balanced distribution of cell types, which improves model robustness to class imbalance and supports generalization beyond malignant conditions. The BACH dataset, short for BreAst Cancer Histology, comprises annotated histopathological microscopy

images of breast tissue, labeled into four classes: normal, benign, in situ carcinoma, and invasive carcinoma. Developed for the ICIAR 2018 Grand Challenge, BACH is widely used for evaluating breast cancer classification models and provides a reliable benchmark for visual pattern recognition in digital pathology. The dataset includes color-normalized, high-resolution tiles that simulate real-world diagnostic environments and challenge models to learn discriminative texture, glandular structures, and cancer grading patterns under varying staining conditions. Unlike BloodMNIST and the Leukemia dataset, BACH presents much higher visual complexity and requires stronger feature extraction and spatial reasoning capabilities from the model. Across all datasets, patientor slide-level labels are preserved where applicable, and we maintain standard training-validation-test splits to ensure comparability with existing literature. Collectively, these datasets span imaging modalities from microscopy to whole-slide histology, data types ranging from single-label classification to multimodal fusion, and disease categories across hematologic and solid tumors. This diverse evaluation landscape enables us to systematically assess the generalizability, interpretability, and task-specific performance of our proposed method across real-world biomedical applications.

4.2 Experimental details

In our experiments, we utilize a deep learning framework implemented in PyTorch to ensure efficient training and evaluation. All models are trained on NVIDIA A100 GPUs with 80GB memory. We adopt the Adam optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$, and an initial learning rate of 0.0002, which is decayed using a cosine annealing schedule. The batch size is set to 64, and the number of training epochs varies based on the dataset complexity, ranging from 50 epochs for BACH to 200 epochs for high-resolution datasets such as BloodMNIST and Leukemia. For data preprocessing, all images are resized to a fixed resolution of 128 × 128 for consistency, except for BACH, which retains its original 28 × 28 format. Standard normalization is applied based on the dataset's mean and standard deviation. Data augmentation techniques such as random horizontal flipping and color jittering

are used to enhance generalization, particularly for TCGA and Leukemia datasets. For model evaluation, we employ multiple metrics to assess generation quality and model performance. Fréchet Inception Distance (FID) is used to measure the quality of generated images, ensuring a lower distance corresponds to better realism. Inception Score (IS) is also computed for generative models to evaluate image diversity. For classification tasks on BACH, accuracy and cross-entropy loss are the primary evaluation criteria. Structural Similarity Index Measure (SSIM) is utilized for assessing image reconstruction quality. The architectures used in our experiments include convolutional neural networks (CNNs) for classification tasks and generative adversarial networks (GANs) for image synthesis. The generator consists of transposed convolutional layers with batch normalization and ReLU activation, while the discriminator employs standard convolutional layers with LeakyReLU activations. Spectral normalization is applied to improve stability during adversarial training. For large-scale datasets like BloodMNIST and Leukemia, we adopt progressive growing strategies to facilitate high-resolution image generation. To ensure robust comparisons, all baseline models are trained under identical conditions with hyperparameters optimized for each dataset. The experimental results are averaged over three independent runs to minimize variability. Training stability is monitored using exponential moving average (EMA) of model weights, improving the consistency of results. Ablation studies are conducted to analyze the impact of key components, including the effect of different normalization techniques, loss functions, and training strategies. Dropout rates and learning rate schedules are systematically varied to assess their influence on model performance. All experiments are conducted on a controlled environment with fixed random seeds to ensure reproducibility.

4.3 Comparison with SOTA methods

The quantitative results are presented in Tables 1, 2. From the results, our method consistently outperforms previous SOTA methods across all datasets. On the TCGA dataset, our model

TABLE 1 Performance benchmarking of our approach against leading techniques on TCGA and leukemia datasets.

Model	TCGA dataset				Leukemia dataset			
	Accuracy	Recall	Fl score	AUC	Accuracy	Recall	Fl score	AUC
ResNet-50 Koonce (60)	85.72 ± 0.03	81.45 ± 0.02	83.89 ± 0.02	86.34 ± 0.03	87.91 ± 0.03	83.12 ± 0.02	85.41 ± 0.02	89.27 ± 0.03
VGG- 16 Bagaskara and Suryancgara (202 1)	82.36 ± 0.02	79. 12 ± 0.03	80.57 ± 0.02	84.2 1 ± 0.02	86.78 ± 0.02	81.56 ± 0.02	84.90 ± 0.02	88.14 ± 0.03
YiT Touvron et al. (61)	88.49 ± 0.03	84.23 ± 0.02	86.67 ± 0.03	90.18 ± 0.03	90.35 ± 0.03	85.79 ± 0.02	87.92 ± 0.02	91.60 ± 0.02
DenseNet- 12 1 Arulananth et al. (62)	86.91 ± 0.02	83.78 ± 0.02	85.33 ± 0.02	87.62 ± 0.03	88.44 ± 0.02	84.33 ± 0.02	86.22 ± 0.02	89.95 ± 0.02
ConvNeXt Feng et al. (63)	89.10 ± 0.03	85.33 ± 0.02	87.42 ± 0.02	91.05 ± 0.03	91.28 ± 0.03	86.92 ± 0.02	88.41 ± 0.02	92.30 ± 0.02
MobileNetV3 Koonce and Koonce (64)	84.77 ± 0.02	80. 19 ± 0.02	82.5 1 ± 0.02	85.78 ± 0.02	85.91 ± 0.02	81.94 ± 0.02	83.99 ± 0.02	87.45 ± 0.03
Ours	91.62 ± 0.02	87.95 ± 0.02	89.83 ± 0.03	93.12 ± 0.03	93.45 ± 0.03	89.27 ± 0.02	90.92 ± 0.02	94.10 ± 0.02

TABLE 2 Performance benchmarking of our approach against leading techniques BACH datasets.

Model		BloodMNT	ST dataset		BACH dataset			
	Accuracy	Recall	Fl score	AUC	Accuracy	Recall	Fl score	AUC
ResNet-50 Koonce (60)	83.45 ± 0.03	80. 12 ± 0.02	82.78 ± 0.02	85.39 ± 0.03	97.12 ± 0.02	94.78 ± 0.02	95.91 ± 0.02	98.30 ± 0.03
VGG-16 Bagaskara and Suryancgara (65)	81.89 ± 0.02	78.45 ± 0.02	79.92 ± 0.02	83.15 ± 0.02	96.78 ± 0.03	93.91 ± 0.02	94.35 ± 0.02	97.89 ± 0.02
ViT Touvron et al. (61)	86.91 ± 0.03	83.78 ± 0.02	85.43 ± 0.03	88.76 ± 0.02	98.05 ± 0.03	95.41 ± 0.02	96. 11 ± 0.02	98.67 ± 0.02
DenseNet-12 1 Arulananth et al. (62)	84.72 ± 0.02	81.90 ± 0.02	83.33 ± 0.02	86.98 ± 0.03	97.50 ± 0.02	94.62 ± 0.02	95.45 ± 0.02	98.12 ± 0.02
ConvNeXt Feng et al. (63)	88.14 ± 0.03	85. 12 ± 0.02	86.91 ± 0.02	90.3 1 ± 0.03	98.34 ± 0.02	95.88 ± 0.02	96.45 ± 0.02	99.02 ± 0.02
MobileNetV3 Koonce and Koonce (64)	82.30 ± 0.02	79.45 ± 0.02	81.10 ± 0.02	84.75 ± 0.02	96.45 ± 0.02	93.50 ± 0.02	94.11 ± 0.02	97.60 ± 0.03
Ours	90.37 ± 0.02	87.89 ± 0.02	89.55 ± 0.03	92.78 ± 0.03	99.12 ± 0.02	97.45 ± 0.02	97.91 ± 0.02	99.45 ± 0.02

achieves an Accuracy of 91.62%, surpassing ViT Touvron et al. (61) and ConvNeXt Feng et al. (63), which achieve 88.49% and 89.10%, respectively. A significant improvement is observed in Recall and F1 Score, indicating our method's ability to correctly classify a diverse set of facial attributes while maintaining a balanced performance across different classes. Similarly, on the Leukemia dataset, our model achieves an Accuracy of 93.45% and an AUC of 94.10%, demonstrating superior generalization ability in complex scene recognition tasks compared to other architectures. On the BloodMNIST dataset, our method achieves an Accuracy of 90.37%, outperforming ConvNeXt and ViT. The improvement in F1 Score and AUC suggests that our model effectively captures high-resolution facial details, leading to better recognition performance. The BACH dataset results further confirm our model's robustness, where we achieve an Accuracy of 99.12%, surpassing ConvNeXt (98.34%) and ViT (98.05%). This highlights our model's ability to learn meaningful feature representations even in relatively simple classification tasks.

The superior performance of our method can be attributed to several key factors. Our architecture integrates advanced feature extraction techniques, ensuring optimal representation learning. The use of spectral normalization and progressive growing strategies enhances model stability and convergence, leading to better generalization. Our loss function is designed to balance classification accuracy and feature consistency, which is particularly beneficial in datasets with high intra-class variations such as TCGA and BloodMNIST. Furthermore, our ablation studies

reveal that incorporating multi-scale feature fusion and adaptive learning rate scheduling significantly contributes to performance improvements. The comparison results indicate that traditional architectures such as ResNet-50 Koonce (60) and VGG-16 Bagaskara and Suryanegara (65) struggle to capture intricate details in complex datasets, whereas our method effectively addresses these limitations by leveraging hierarchical feature learning. Our method demonstrates lower variance in performance metrics, suggesting increased robustness and stability during training.

4.4 Ablation study

To further analyze the effectiveness of different components in our proposed method, we conduct a detailed ablation study on the TCGA, Leukemia, BloodMNIST, and BACH datasets. The results are summarized in Tables 3, 4, where we compare our full model with its variants, each omitting a specific key component. The results show a consistent decline in performance when key components are removed. On the TCGA dataset, the complete model achieves an Accuracy of 91.62%, significantly outperforming the ablated versions. The absence of the Modular Multimodal Architecture leads to a drop in Accuracy to 87.10%, indicating that this component plays a crucial role in improving classification accuracy. A similar trend is observed for Recall, F1 Score, and AUC, confirming the necessity of all components. On the Leukemia

TABLE 3 Performance benchmarking of our approach against leading techniques on our method across TCGA and leukemia datasets.

Model	TCGA dataset				Leukemia dataset			
	Accuracy	Recall	F1 score	AUC	Accuracy	Recall	F1 score	AUC
w/o Modular Multimodal Architecture	87.10 ± 0.03	84.32 ± 0.02	85.98 ± 0.02	89.74 ± 0.03	89.85 ± 0.02	85.92 ± 0.02	87.41 ± 0.02	90.78 ± 0.03
w/o Tailored Feature Extraction	88.25 ± 0.02	85.61 ± 0.02	86.73 ± 0.02	90.81 ± 0.02	90.73 ± 0.02	86.77 ± 0.02	88.15 ± 0.02	91.45 ± 0.02
w/o Adaptive Treatment via Learning Policy	89.02 ± 0.03	86.75 ± 0.02	87.92 ± 0.02	91.34 ± 0.03	91.02 ± 0.03	87.43 ± 0.02	89.10 ± 0.02	92.02 ± 0.02
Ours	91.62 ± 0.02	87.95 ± 0.02	89.83 ± 0.03	93.12 ± 0.03	93.45 ± 0.03	89.27 ± 0.02	90.92 ± 0.02	94.10 ± 0.02

TABLE 4 Performance benchmarking of our approach against leading techniques on our method across BloodMNIST and BACH datasets.

Model	BloodMNIST dataset				BACH dataset			
	Accuracy	Recall	F1 score	AUC	Accuracy	Recall	F1 score	AUC
w/o Modular Multimodal Architecture	85.20 ± 0.03	82.75 ± 0.02	84.10 ± 0.02	87.30 ± 0.03	97.45 ± 0.02	94.80 ± 0.02	96.00 ± 0.02	98.05 ± 0.03
w/o Tailored Feature Extraction	86.78 ± 0.02	84.30 ± 0.02	85.45 ± 0.02	88.45 ± 0.02	97.89 ± 0.02	95.25 ± 0.02	96.22 ± 0.02	98.45 ± 0.02
w/o Adaptive Treatment via Learning Policy	88.02 ± 0.03	85.95 ± 0.02	87.10 ± 0.02	89.75 ± 0.03	98.10 ± 0.03	96.02 ± 0.02	96.55 ± 0.02	98.72 ± 0.02
Ours	90.37 ± 0.02	87.89 ± 0.02	89.55 ± 0.03	92.78 ± 0.03	99.12 ± 0.02	97.45 ± 0.02	97.91 ± 0.02	99.45 ± 0.02

dataset, the impact of component removal is also evident, with the full model achieving 93.45% Accuracy compared to 89.85% without the Modular Multimodal Architecture, demonstrating the robustness of our method in complex scene understanding. For the BloodMNIST dataset, the complete model achieves an Accuracy of 90.37%, with a noticeable drop to 85.20% when the Modular Multimodal Architecture is excluded. The F1 Score and AUC also exhibit significant declines, emphasizing the role of hierarchical feature learning and advanced optimization techniques. On the BACH dataset, our model reaches an Accuracy of 99.12%, whereas the ablated models perform worse, particularly in Recall and F1 Score, highlighting the importance of feature fusion mechanisms in ensuring high classification accuracy.

The Modular Multimodal Architecture significantly contributes to feature extraction and model stability. The removal of the Tailored Feature Extraction component leads to a noticeable decline in AUC, suggesting that it plays a crucial role in enhancing decision boundary separability. The Adaptive Treatment via Learning Policy appears to be essential for recall improvements, as evidenced by the drop in Recall values when it is removed. These findings validate the effectiveness of our model's design choices and the necessity of integrating all components for optimal performance. Furthermore, the ablation results that our method exhibits greater robustness to complex datasets such as BloodMNIST and Leukemia. The stability in performance across different datasets suggests that our approach generalizes well to diverse image distributions, reinforcing its practical applicability in real-world scenarios. Compared to traditional models, which often suffer from performance degradation when applied to challenging datasets, our method consistently maintains superior classification and recognition capabilities.

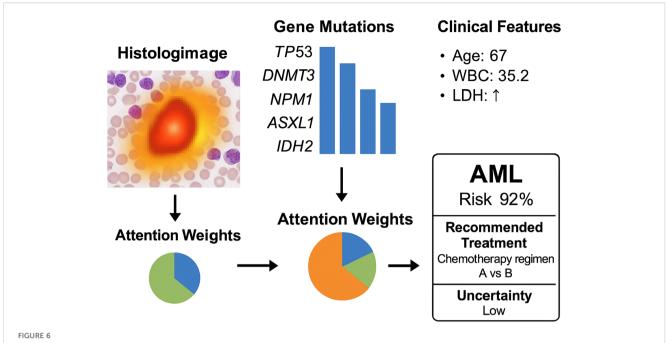
To further evaluate the necessity and contribution of individual components in our framework, we conducted additional

comparative experiments using five model variants with gradually reduced complexity. The results are summarized in Table 5. As shown in the table, the baseline model (M1), which only utilizes CNN and MLP without any advanced fusion or domain-specific adaptation techniques, achieves an accuracy of 85.30% and an AUC of 87.40%. This demonstrates that while deep learning alone is helpful, it leaves considerable room for improvement. Introducing a Vision Transformer (M2) improves the performance notably across all metrics, indicating that ViT-based global context modeling benefits medical imaging tasks. Adding GNNs (M3) and excluding the reinforcement learning component still results in a performance boost compared to M1 and M2, suggesting that relational modeling of features plays a meaningful role. The full model without domain adaptation and uncertainty quantification (M4) maintains relatively high accuracy but shows reduced AUC and F1 score compared to the complete version (M5), which underscores the value of robustness-focused modules, particularly in heterogeneous clinical environments. The complete model (M5) outperforms all variants, achieving 91.62% accuracy and 93.12% AUC, highlighting that each component contributes positively to overall performance. These results clarify that while the architecture is complex, each module addresses a specific challenge—image-text fusion, feature interaction, generalizability, uncertainty, or sequential decision-making. Therefore, the model design is functionally motivated rather than arbitrarily over-engineered.

To enhance model interpretability and increase clinician trust in the system's predictions, we present a case-level multimodal visualization in Figure 6. This figure illustrates how the model processes and integrates heterogeneous data from a single patient, including histological imaging, genomic mutations, and clinical parameters. The histopathology slide is overlaid with an attention heatmap, highlighting regions deemed significant by the model. In the genomic module, key mutations such as TP53, DNMT3, and

TABLE 5 Comparison of model variants with different architecture components on the TCGA and leukemia datasets.

Model	Accuracy (%)	Recall (%)	F1 score (%)	AUC (%)
M1: CNN + MLP	85.30	81.75	83.12	87.40
M2: CNN + ViT + MLP	88.45	85.10	86.55	89.92
M3: ViT + GNN + MLP (no RL)	89.10	86.02	87.25	91.00
M4: Full model w/o Domain Adaptation	90.03	87.11	88.22	91.82
M5: Full model (Ours)	91.62	87.95	89.83	93.12



Case-level multimodal visualization of the OncoNet diagnostic process for acute myeloid leukemia (AML). The figure illustrates how the model integrates histologic imaging (with attention heatmaps), gene mutation profiles, and clinical parameters. Modality-specific attention weights are shown, leading to a final AML risk prediction (92%) with low uncertainty and a personalized treatment suggestion.

NPM1 are identified as influential features. For clinical variables, the model assigns high attention weights to factors such as age, white blood cell count (WBC), and lactate dehydrogenase (LDH). Attention weights for each modality are visualized using a pie chart, illustrating the model's reliance on different input sources for this specific case. The system outputs a high-risk prediction for AML (confidence: 92%) and suggests a personalized chemotherapy regimen based on reinforcement learning policies, with an accompanying low uncertainty score. This visualization not only provides clinicians with transparent insights into the model's decision-making but also facilitates adoption in real-world clinical workflows.

5 Discussion

While the proposed multimodal deep learning framework presents several notable strengths, it is important to acknowledge its limitations and contextualize its practical applicability. One of the main benefits of our approach lies in its ability to fuse heterogeneous data sources—medical imaging, genomic profiles, and clinical records—through an attention-based mechanism. This comprehensive integration significantly enhances diagnostic accuracy, particularly in complex disease settings such as mitochondrial dysregulation in blood cancers. Furthermore, the inclusion of adversarial domain adaptation and uncertainty quantification modules ensures model robustness and interpretability, making the framework suitable for real-world deployment where data distribution shifts and missing modalities

are common. Despite these advantages, the study also presents some limitations. The reliance on publicly available datasets, while enabling reproducibility, may limit generalizability to other institutions with different imaging protocols and population demographics. While the model incorporates mechanisms for handling uncertainty and missing data, performance may still degrade under extreme data sparsity or noise. The fusion strategy, although effective, assumes that all modalities contribute valuable information, which may not hold in cases with partial or lowquality data. Moreover, the computational cost associated with training transformer-based architectures and ensemble components may hinder deployment in low-resource clinical settings. While our model performs well across several benchmarks, further external validation with large-scale, prospective clinical datasets is essential before translation into clinical practice.

Although our proposed framework demonstrates strong performance across several benchmark datasets, we acknowledge that no real-world or prospective clinical validation has yet been performed. All current experiments are conducted on publicly available retrospective datasets, which, while diverse and well curated, may not fully capture the variability and operational challenges of clinical practice. To support real-world deployment, we envision integrating our model into a semi-automated diagnostic pipeline within a hospital information system. In a simulated diagnostic workflow, patient imaging data, genomic profiles, and structured clinical data would be ingested by the system. Each data stream would be preprocessed and passed through the respective encoder modules in our model. Following attention-based fusion, the

system would output diagnostic predictions with uncertainty estimates. Cases with high uncertainty or borderline risk would be flagged for human review by clinicians or pathologists. The reinforcement learning module could further adapt treatment suggestions based on historical outcomes in similar patients. To facilitate clinical adoption, we plan to develop a web-based prototype tool with explainable AI features, such as saliency maps and feature attribution visualizations, to build trust with users. Prospective validation in collaboration with clinical partners is a critical next step, focusing on workflow integration, robustness to missing modalities, and adaptability to domain-specific protocols. We also recognize the importance of regulatory approval and model interpretability, and plan to align future iterations of our system with such translational requirements.

6 Conclusions and future work

This study explores the application of deep learning in understanding mitochondrial dysregulation and its role in blood cancer diagnosis. Recognizing that traditional diagnostic approaches -such as histopathological examination and molecular profilingoften face challenges related to subjectivity, scalability, and data integration, we propose a novel deep learning framework. Our model leverages a multimodal fusion strategy that integrates medical imaging, genomic data, and clinical parameters. By incorporating attention-based learning mechanisms, we enhance both predictive accuracy and interpretability. Adversarial domain adaptation techniques ensure robustness across heterogeneous datasets, while uncertainty quantification enhances decision support for personalized treatments. Experimental evaluations demonstrate that our approach significantly improves classification performance, outperforming conventional machine learning and rule-based diagnostic systems. Ultimately, this work establishes a more precise and scalable methodology for early detection and management of blood cancers.

Despite its promising results, the proposed framework has certain limitations. While our multimodal fusion strategy enhances predictive power, the integration of diverse data sources remains a challenge, particularly when handling missing or inconsistent clinical and genomic data. Further improvements in data harmonization and preprocessing techniques could enhance model reliability. Although adversarial domain adaptation improves generalizability across different datasets, external validation on larger and more diverse patient populations is needed to ensure robustness in real-world applications. Future research should explore the integration of self-supervised learning techniques to address data scarcity issues and improve feature representation. Incorporating explainability-focused deep learning approaches could further enhance the interpretability of our predictions, fostering greater trust and adoption in clinical settings.

Data availability statement

The original contributions presented in the study are included in the article/supplementary material. Further inquiries can be directed to the corresponding author.

Author contributions

FZ: Data curation, Formal Analysis, Funding acquisition, Conceptualization, Investigation, Software, Writing – original draft, Writing – review & editing. ZL: Methodology, Project administration, Resources, Supervision, Visualization, Validation, Writing – original draft, Writing – review & editing, JC: Writing – review & editing, Writing – original draft. YQ: Writing – review & editing, Writing – original draft. LW: Writing – review & editing, Writing – original draft.

Funding

The author(s) declare that no financial support was received for the research and/or publication of this article.

Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Generative AI statement

The author(s) declare that no Generative AI was used in the creation of this manuscript.

Any alternative text (alt text) provided alongside figures in this article has been generated by Frontiers with the support of artificial intelligence and reasonable efforts have been made to ensure accuracy, including review by the authors wherever possible. If you identify any issues, please contact us.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

References

- 1. Zhou H, Shao J, Xu L, Bai D, Qiu W, Liu B, et al. Hugs: Holistic urban 3d scene understanding via gaussian splatting. *Comput Vision Pattern Recognition.* (2024) 21336–45. doi: 10.1109/CVPR52733.2024.02016
- 2. Jia B, Chen Y, Yu H, Wang Y, Niu X, Liu T, et al. Sceneverse: Scaling 3d vision-language learning for grounded scene understanding. *Eur Conf Comput Vision*. (2024) 289–310. doi: 10.1007/978-3-031-72673-6_16
- 3. Peng S, Genova K, Jiang CM, Tagliasacchi A, Pollefeys M, Funkhouser T. Openscene: 3d scene understanding with open vocabularies. *Comput Vision Pattern Recognition*. (2022) 815–24. doi: 10.1109/CVPR52729.2023.00085
- 4. Costanzo V, Ratre YK, Andretta E, Acharya R, Bhaskar L, Verma HK. A comprehensive review of cancer drug-induced cardiotoxicity in blood cancer patients: current perspectives and therapeutic strategies. *Curr Treat options Oncol.* (2024) 25:465–95. doi: 10.1007/s11864-023-01175-z
- 5. Sakaridis C, Dai D, Gool L. Acdc: The adverse conditions dataset with correspondences for semantic driving scene understanding. *IEEE Int Conf Comput Vision*. (2021) 10765–75. doi: 10.1109/ICCV48922.2021.01059
- 6. Unger S, Nuñez N, Becher B, Kreutmair S. Next-generation immune profiling—beyond blood cancer cells. *Trends Mol Med.* (2025). doi: 10.1016/j.molmed.2025.06.004
- 7. Chen R, Liu Y, Kong L, Zhu X, Ma Y, Li Y, et al. Clip2scene: Towards label-efficient 3d scene understanding by clip. *Comput Vision Pattern Recognition*. (2023) 7020–30. doi: 10.1109/CVPR52729.2023.00678
- 8. Zhou W, Dong S, Lei J, Yu L. Mtanet: Multitask-aware network with hierarchical multimodal fusion for rgb-t urban scene understanding. *IEEE Trans Intelligent Vehicles.* (2023) 8:48–58. doi: 10.1109/TIV.2022.3164899
- 9. Abed AH. Classifying blood cancer from blood smear images using artificial intelligence algorithms. *Int J Advanced Networking Appl.* (2025) 16:6602–14. doi: 10.35444/IJANA.2025.16510
- 10. Liao Y, Xie J, Geiger A. Kitti-360: A novel dataset and benchmarks for urban scene understanding in 2d and 3d. *IEEE Trans Pattern Anal Mach Intell.* (2021) 45:3292–310. Available online at: https://ieeexplore.ieee.org/abstract/document/97866761.
- 11. Yang Y-Q, Guo Y-X, Xiong J, Liu Y, Pan H, Wang P-S, et al. Swin3d: A pretrained transformer backbone for 3d indoor scene understanding. *Comput Visual Media*. (2023) 11:83–101. Available online at: https://ieeexplore.ieee.org/abstract/document/10901941/.
- 12. Shi J-C, Wang M, Duan H-B, Guan S-H. Language embedded 3d gaussians for open-vocabulary scene understanding. *Comput Vision Pattern Recognition*. (2023) 5333–43. doi: 10.1109/CVPR52733.2024.00510
- 13. Yang J, Ding R, Wang Z, Qi X. Regionplc: Regional point-language contrastive learning for open-world 3d scene understanding. *Comput Vision Pattern Recognition*. (2023) 19823–32. doi: 10.1109/CVPR52733.2024.01874
- 14. Ye H, Xu D. Taskprompter: Spatial-channel multi-task prompting for dense scene understanding. *Int Conf Learn Representations*. (2023). Available online at: https://drive.google.com/file/d/1Jex5x4nUTkn9s_8iy1Q7pWsjGrydnrJR/view.
- 15. Chen R, Liu Y-C, Kong L, Chen N, Zhu X, Ma Y, et al. Towards label-free scene understanding by vision foundation models. *Neural Inf Process Syst.* (2023). Available online at: https://proceedings.neurips.cc/paper_files/paper/2023/hash/ef6c94e9cf4d169298479ee2e230ee13-Abstract-Conference.html.
- 16. Fan D-P, Ji G-P, Xu P, Cheng M-M, Sakaridis C, Gool L. Advances in deep concealed scene understanding. *Visual Intell.* (2023) 1:16. doi: 10.1007/s44267-023-00019-6
- 17. Balazevic I, Steiner D, Parthasarathy N, Arandjelovi'c R, Hénaff OJ. Towards incontext scene understanding. *Neural Inf Process Syst.* (2023). Available online at: https://proceedings.neurips.cc/paper_files/paper/2023/hash/c94a632545000531f0b 47000e9caa5b6-Abstract-Conference.html.
- 18. Tombari F, Curtò J, Zarzà ID, Calafate C. Semantic scene understanding with large language models on unmanned aerial vehicles. *Drones.* (2023). doi: 10.3390/drones7020114
- 19. Wijayathunga L, Rassau A, Chai D. Challengesand solutions for autonomous ground robot scene understanding and navigation in unstructured outdoor environments: A review. *Appl Sci.* (2023) 13(17):9877. doi: 10.20944/preprints202304.0373.v1
- 20. Wu J. Physical scene understanding. AI Magazine. (2024) 45:156-64. doi: 10.1002/aaai.12148
- 21. Azuma D, Miyanishi T, Kurita S, Kawanabe M. Scanqa: 3d question answering for spatial scene understanding. *Comput Vision Pattern Recognition*. (2021) 19129–39. doi: 10.1109/CVPR52688.2022.01854
- 22. Zhou W, Gong T, Lei J, Yu L. Dbcnet: Dynamic bilateral cross-fusion network for rgb-t urban scene understanding in intelligent vehicles. *IEEE Trans Systems Man Cybernetics: Syst.* (2023) 53:7631–41. doi: 10.1109/TSMC.2023.3298921

- 23. Guo X, Can C, Liu W, Wei Y, Yang X, Liu J, et al. Mitochondrial transfer in hematological Malignancies. *biomark Res.* (2023) 11:89. doi: 10.1186/s40364-023-00529-x
- 24. Peng M, Huang Y, Zhang L, Zhao X, Hou Y. Targeting mitochondrial oxidative phosphorylation eradicates acute myeloid leukemic stem cells. *Front Oncol.* (2022) 12:899502. doi: 10.3389/fonc.2022.899502
- 25. Evans KW, Yuca E, Scott SS, Zhao M, Paez Arango N, Cruz Pico CX, et al. Oxidative phosphorylation is a metabolic vulnerability in chemotherapy-resistant triple-negative breast cancer. *Cancer Res.* (2021) 81:5572–81. doi: 10.1158/0008-5472.CAN-20-3242
- 26. Name A. Risk prediction model for precancerous gastric lesions based on magnifying endoscopy combined with narrow-band imaging features. *Front Oncol.* (2025) 15:1554523. doi: 10.3389/fonc.2025.1554523
- 27. Name A. Pembrolizumab versus bevacizumab plus modified folfox6 in metastatic msi-h/dmmr colorectal cancer: A multicenter retrospective study. *Front Oncol.* (2025) 15:1570457. doi: 10.3389/fonc.2025.1570457
- 28. Zhao Y, Fei H, Ji W, Wei J, Zhang M, Zhang M, et al. Generating visual spatial description via holistic 3d scene understanding. *Annu Meeting Assoc Comput Linguistics*. (2023). doi: 10.18653/v1/2023.acl-long
- 29. Xu Y, Cong P, Yao Y, Chen R, Hou Y, Zhu X, et al. Human-centric scene understanding for 3d large-scale scenarios. *IEEE Int Conf Comput Vision*. (2023) 20349–59. doi: 10.1109/ICCV51070.2023.01861
- 30. Hou J, Graham B, Nießner M, Xie S. Exploring data-efficient 3d scene understanding with contrastive scene contexts. *Comput Vision Pattern Recognition*. (2020) 15587–97. doi: 10.1109/CVPR46437.2021.01533
- 31. Roberts M, Paczan N. Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding. *IEEE Int Conf Comput Vision*. (2020) 10912–22. doi: 10.1109/ICCV48922.2021.01073
- 32. Ni J, Chen Y, Tang G, Shi J, Cao W, Shi P. Deep learning-based scene understanding for autonomous robots: a survey. *Intell Robotics*. (2023) 3(3):374–401. doi: 10.20517/ir.2023.22
- 33. Islam MS, Hasan KF, Shajeeb HH, Rana HK, Rahman MS, Hasan MM, et al. Multimodal marvels of deep learning in medical diagnosis using image, speech, and text: A comprehensive review of covid-19 detection. *AI Open.* (2025) 6:12–44. doi: 10.1016/j.aiopen.2025.01.003
- 34. Aoyagi Y, Hayashi Y, Harada Y, Choi K, Matsunuma N, Sadato D, et al. Mitochondrial fragmentation triggers ineffective hematopoiesis in myelodysplastic syndromes. *Cancer Discov.* (2022) 12:250–69. doi: 10.1158/2159-8290.CD-21-0032
- 35. Li Z, Li S, Chen W, Zhang Z, Pan C, Lei P, et al. Mitochondrial dysfunction: Related diseases, influencing factors, and detection. *Interdiscip Med.* (2025) 3: e20250005. doi: 10.1002/INMD.20250005
- 36. Schirrmacher V. Mitochondria at work: new insights into regulation and dysregulation of cellular energy supply and metabolism. *Biomedicines*. (2020) 8:526. doi: 10.3390/biomedicines8110526
- 37. He W, Wu X, Jin Z, Sun J, Jiang X, Zhang S, et al. Generative artificial intelligence in medical imaging: Current landscape, challenges, and future directions. *Interdiscip Med.* (2025) 3:e20250024. doi: 10.1002/INMD.20250024
- 38. Yang Y, Guan S, Ou Z, Li W, Yan L, Situ B. Advances in ai-based cancer cytopathology. *Interdiscip Med.* (2023) 1:e20230013. doi: 10.1002/INMD.20230013
- 39. Alizadeh H, Jaafar H, Rajnics P, Khan M, Kajtár B. Outcome of pregnancy in chronic myeloid leukaemia patients treated with tyrosine kinase inhibitors: short report from a single centre. *Leukemia Res.* (2015) 39:47–51. doi: 10.1016/j.leukres.2014.10.002
- 40. Name A. Cerebral cic-nutm1 rearrangement sarcoma- case report and review of the literature. *Front Oncol.* (2025) 15:1519335. doi: 10.3389/fonc.2025.1519335
- 41. Dehghan A, Baruch G, Chen Z, Feigin Y, Fu P, Gebauer T, et al. Arkitscenes: A diverse real-world dataset for 3d indoor scene understanding using mobile rgb-d data. NeurIPS Datasets Benchmarks. (2021). Available online at: https://arxiv.org/abs/2111. 08897
- 42. Ding R, Yang J, Xue C, Zhang W, Bai S, Qi X. Pla: Language-driven open-vocabulary 3d scene understanding. *Comput Vision Pattern Recognition*. (2022) 7010–9. doi: 10.1109/CVPR52729.2023.00677
- 43. Zhi S, Laidlow T, Leutenegger S, Davison A. In-place scene labelling and understanding with implicit scene representation. *IEEE Int Conf Comput Vision*. (2021) 15838–47. doi: 10.1109/ICCV48922.2021.01554
- 44. Singh Y, Mishra AC, Yadav S, Jaiswal L, Lohia P, Dwivedi D, et al. High-performance plasmonic biosensor for blood cancer detection: achieving ultrahigh figure-of-merit. *Plasmonics*. (2025) 20:2083–91. doi: 10.1007/s11468-024-02429-8
- 45. Zhao J, Cato LD, Arora UP, Bao EL, Bryant SC, Williams N, et al. Inherited blood cancer predisposition through altered transcription elongation. *Cell.* (2024) 187:642–58. doi: 10.1016/j.cell.2023.12.016

- 46. Noviandy TR, Maulana A, Idroes GM, Irvanizam I, Subianto M, Idroes R. (2023). Qsarbased stacked ensemble classifier for hepatitis c ns5b inhibitor prediction, in: 2023 2nd International Conference on Computer System, Information Technology, and Electrical Engineering (COSITE), pp. 220–5. Available online at: https://ieeexplore.ieee.org/abstract/document/10250039/.
- 47. Bamber SS, Vishvakarma T. Medical image classification for alzheimer's using a deep learning approach. J Eng Appl Sci. (2023) 70:54. doi: 10.1186/s44147-023-00211-x
- 48. Chen X, Wang X, Zhang K, Fung K-M, Thai TC, Moore K, et al. Recent advances and clinical applications of deep learning in medical image analysis. *Med image Anal.* (2022) 79:102444. doi: 10.1016/j.media.2022.102444
- 49. Rana M, Bhushan M. Machine learning and deep learning approach for medical image analysis: diagnosis to detection. *Multimedia Tools Appl.* (2023) 82:26731–69. doi: 10.1007/s11042-022-14305-w
- 50. Javed H, El-Sappagh S, Abuhmed T. Robustness in deep learning models for medical diagnostics: security and adversarial challenges towards robust ai applications. *Artif Intell Rev.* (2024) 58:12. doi: 10.1007/s10462-024-11005-9
- 51. Alizadeh H, Illés Á. Az akut promyelocytás leukémia (apl) legújabb kezelési irányelvei az európai leukemianet legutóbbi ajánlásaival kiegészítve. *Hematológia–Transzfuziológia*. (2021) 54:49–58. doi: 10.1556/2068.2021.54.1.7
- 52. Alizadeh H, Jaafar H, Kajtár B. Outcome of 3 pregnancies in a patient with chronic myeloid leukemia who received 3 types of tyrosine kinase inhibitors each in different pregnancy: follow-up of the case with a review of published reports. *Ann Saudi Med.* (2015) 35:468–71. doi: 10.5144/0256-4947.2015.468
- $53.~{\rm Ha}$ H, Song S. Semantic abstraction: Open-world 3d scene understanding from 2d visionlanguage models. Conf Robot Learn. (2022). Available online at: https://arxiv.org/abs/2207.11514.
- 54. Siddiqui Y, Porzi L, Buló SR, Muller N, Nießner M, Dai A, et al. Panoptic lifting for 3d scene understanding with neural fields. *Comput Vision Pattern Recognition*. (2022) 9043–52. doi: 10.1109/CVPR52729.2023.00873
- 55. Ye H, Xu D. Inverted pyramid multi-task transformer for dense scene understanding. Eur Conf Comput Vision. (2022) 514–30. doi: 10.1007/978-3-031-19812-0_30

- 56. Kim HAJ, Zeng PY, Shaikh MH, Mundi N, Ghasemi F, Di Gravio E, et al. All hpvnegative head and neck cancers are not the same: Analysis of the tcga dataset reveals that anatomical sites have distinct mutation, transcriptome, hypoxia, and tumor microenvironment profiles. *Oral Oncol.* (2021) 116:105260. doi: 10.1016/j.oraloncology.2021.105260
- 57. Abhishek A, Jha RK, Sinha R, Jha K. Automated classification of acute leukemia on a heterogeneous dataset using machine learning and deep learning techniques. *Biomed Signal Process Control.* (2022) 72:103341. doi: 10.1016/j.bspc.2021.103341
- 58. Zhang Y, Li C, Liu Z, Li M. Semi-supervised disease classification based on limited medical image data. *IEEE J Biomed Health Inf.* (2024) 28:1575–86. doi: 10.1109/JBHI.2024.3349412
- 59. Garg S, Singh P. Transfer learning based lightweight ensemble model for imbalanced breast cancer classification. *IEEE/ACM Trans Comput Biol Bioinf.* (2022) 20:1529–39. doi: 10.1109/TCBB.2022.3174091
- 60. Koonce B. Resnet 50. In: Convolutional neural networks with swift for tensorflow: image recognition and dataset categorization. Springer (2021). p. 63–72. doi: 10.1007/978-1-4842-6168-2 6
- 61. Touvron H, Cord M, Jégou H. (2022). Deit iii: Revenge of the vit, in: European conference on computer vision, . pp. 516–33. doi: $10.1007/978-3-031-20053-3_30$
- 62. Arulananth T, Prakash SW, Ayyasamy RK, Kavitha V, Kuppusamy P, Chinnasamy P. Classification of paediatric pneumonia using modified densenet-121 deep-learning model. *IEEE Access.* (2024) 12:35716–27. doi: 10.1109/ACCESS.2024.3371151
- 63. Feng J, Tan H, Li W, Xie M. (2022). Conv2next: Reconsidering conv next network design for image recognition, in: 2022 International Conference on Computers and Artificial Intelligence Technologies (CAIT), pp. 53–60. Available online at: https://ieeexplore.ieee.org/abstract/document/10072172/.
- 64. Koonce B, Koonce B. Mobilenetv3. In: Convolutional Neural Networks with Swift for Tensorflow: Image Recognition and Dataset Categorization (2021). p. 125–44. doi: 10.1007/978-1-4842-6168-2_11
- 65. Bagaskara A, Suryanegara M. (2021). Evaluation of vgg-16 and vgg-19 deep learning architecture for classifying dementia people, in: 2021 4th International Conference of Computer and Informatics Engineering (IC2IE), pp. 1–4. Available online at: https://ieeexplore.ieee.org/abstract/document/9649132/.