Check for updates

# Anomaly detection in medical via multimodal foundation models

Zhenyou Tang[1], Zhong Tang[2]* and Jing Wu[3]

[1]Institute of Collaborative Innovation, University of Macau, Macau, China, [2]College of Humanities and Social Sciences, Guangxi Medical University, Guangxi, Nanning, China, [3]Xi'an University of Architecture and Technology, Xi'an, China

**Introduction:** Recent advances in artificial intelligence have created opportunities for medical anomaly detection through multimodal learning frameworks. However, traditional systems struggle to capture the complex temporal and semantic relationships in clinical data, limiting generalization and interpretability in real-world settings.

**Methods:** To address these challenges, we propose a novel framework that integrates symbolic representations, a graph-based neural model (PathoGraph), and a knowledge-guided refinement strategy. The approach leverages structured clinical records, temporally evolving symptom graphs, and medical ontologies to build semantically interpretable latent spaces. Our method enhances model robustness under sparse supervision and distributional shifts.

**Results:** Extensive experiments across electronic health records and diagnostic datasets show that our model outperforms existing baselines in detecting rare comorbidity patterns and abnormal treatment responses.

**Discussion:** Additionally, it improves interpretability and trustworthiness, which are critical for clinical deployment. By aligning domain knowledge with multimodal AI, our work contributes a generalizable and explainable solution to healthcare anomaly detection.

## 1 Introduction

Anomaly detection in the medical domain is crucial for ensuring early diagnosis and timely intervention, thereby significantly improving patient outcomes. With the proliferation of digital healthcare data-from radiology images to clinical notes-there exists an unprecedented opportunity to enhance anomaly detection using computational techniques (Liu Z. et al., 2023). Traditional approaches often struggle to integrate diverse data types or generalize across clinical settings. Not only are many conventional models modality-specific, limiting their capacity to capture complex cross-modal patterns, but they also require extensive manual feature engineering (Roth et al., 2021). Recent advances in artificial intelligence, particularly in the domain of multimodal foundation models, offer a promising path forward. These models can seamlessly integrate and reason over heterogeneous data sources such as text, images, and signals, providing a unified framework for medical anomaly detection (Deng and Li, 2022). Moreover, they can leverage pretraining on vast datasets to generalize across tasks with minimal supervision, which is particularly valuable in medical contexts where annotated data are scarce (Zou et al., 2022).

Initial efforts to detect anomalies in medical settings began with systems that relied on predefined rules and clinical coding hierarchies to identify deviations from normal health states (Tuli et al., 2022). These implementations were highly interpretable and mirrored human expert reasoning but struggled to remain effective when faced with incomplete records, noisy measurements, or evolving clinical practices (Li et al., 2021). They lacked the flexibility to incorporate new forms of data, such as imaging or biosignals, and often failed to scale across medical specialties. For instance, systems calibrated for cardiovascular monitoring had limited utility when applied to neuroimaging diagnostics (Zavrtanik et al., 2021).

To expand applicability and reduce dependence on hand-crafted knowledge, researchers began leveraging statistical learning techniques that could infer patterns from empirical examples (Liu J. et al., 2023). Algorithms were trained to discriminate between normal and abnormal health indicators using structured datasets such as laboratory results or physiological waveforms (Deng and Hooi, 2021). This marked a step toward greater adaptability, yet these systems still faced critical limitations (Wang D. et al., 2023). Designing input features remained a manual, expertise-intensive task, and most models operated in isolation on single-modality data, leaving valuable cross-domain correlations untapped. As a result, their ability to detect subtle or complex clinical anomalies remained constrained (You et al., 2022).

The rapid growth in computational power and data availability eventually enabled a shift toward more expressive models capable of learning directly from raw inputs (Gudovskiy et al., 2021). Neural networks-especially convolutional, recurrent, and transformer-based architectures-ushered in new possibilities for identifying anomalies across a range of medical domains, from radiological imaging to genomic sequences (Tian et al., 2021). More recently, the emergence of large-scale multimodal models has allowed for the joint analysis of text, images, and signals under a unified computational framework. These models are pretrained on diverse medical corpora and refined on specific tasks, offering superior generalization with minimal supervision (Bergmann et al., 2021). While they pose challenges in terms of interpretability and computational demands, their capacity to capture complex interdependencies across data types positions them as the most promising approach for future medical anomaly detection systems (Liu et al., 2021).

Nevertheless, many current models still fall short in robustness, semantic consistency, and alignment with clinical reasoning processes. To overcome these limitations, we propose a novel anomaly detection framework built upon a multimodal foundation model architecture, tailored to the complexities of medical data.

Our system integrates three synergistic modules: (1) a mathematically formalized symbolic abstraction of multimodal clinical records; (2) PathoGraph, a graph-based neural model that constructs a temporally-evolving, symptom-centric latent space for structured disentanglement; and (3) Knowledge-Guided Refinement (KGR), a strategic overlay that embeds domain ontologies such as SNOMED CT and ICD-10 into the learning pipeline via differentiable constraints and uncertainty-aware attention mechanisms.

This integrative design enhances detection performance while ensuring semantic interpretability and clinical plausibility. The proposed framework demonstrates superior results across multiple real-world diagnostic datasets, successfully identifying rare and complex anomalies under weak supervision, while maintaining alignment with symbolic medical knowledge. The proposed approach offers a range of significant benefits that set it apart from conventional methods.

- We introduce a novel cross-modal attention module that dynamically integrates visual, textual, and physiological features, offering a unified and context-aware representation for anomaly detection.
- Our model excels in multi-scenario deployment, demonstrating high efficiency and generalization across diverse clinical tasks, from radiology to pathology.
- Experiments on multiple public and proprietary datasets show significant improvements in detection precision and recall, outperforming state-of-the-art baselines in multimodal anomaly detection.

# 2 Related work

## 2.1 Multimodal learning in medicine

Multimodal learning has emerged as a critical paradigm in medical artificial intelligence, enabling the integration of heterogeneous data sources such as medical images, electronic health records (EHRs), clinical notes, and genomic data (Liu et al., 2025). This integration allows for richer representations that facilitate improved diagnostic accuracy and patient outcome predictions (Tang et al., 2022). In recent years, large-scale multimodal foundation models have demonstrated an exceptional ability to encode cross-modal information through unified architectures, such as transformers, that jointly learn from text and images (Bayane et al., 2025). For instance, models like CLIP and MedCLIP adapt the contrastive learning framework to align visual and textual modalities in the medical domain (Mishra et al., 2021). These methods leverage large-scale, weakly labeled datasets to learn generalizable representations without extensive annotation. In the medical context, multimodal models have been applied to tasks including radiology report generation, disease classification, and decision support (Zhou et al., 2023). Such models have shown the capacity to capture nuanced correlations across modalities, such as linking radiological patterns with specific terminologies in textual reports. Transfer learning and domain adaptation strategies are often employed to enhance model robustness across different medical subdomains or imaging modalities (Jiang et al., 2023). Moreover, recent advancements have focused on designing unified pretraining objectives that incorporate both contrastive and generative tasks, leading to more comprehensive embeddings (Tien et al., 2023). One notable challenge is modality-specific noise and missing data. Medical data is often incomplete or irregularly sampled across patients. Techniques such as modality dropout, modality-aware fusion mechanisms, and imputation with attention have been introduced to address these issues (Makrogiannis et al., 2021). Despite promising results, evaluating
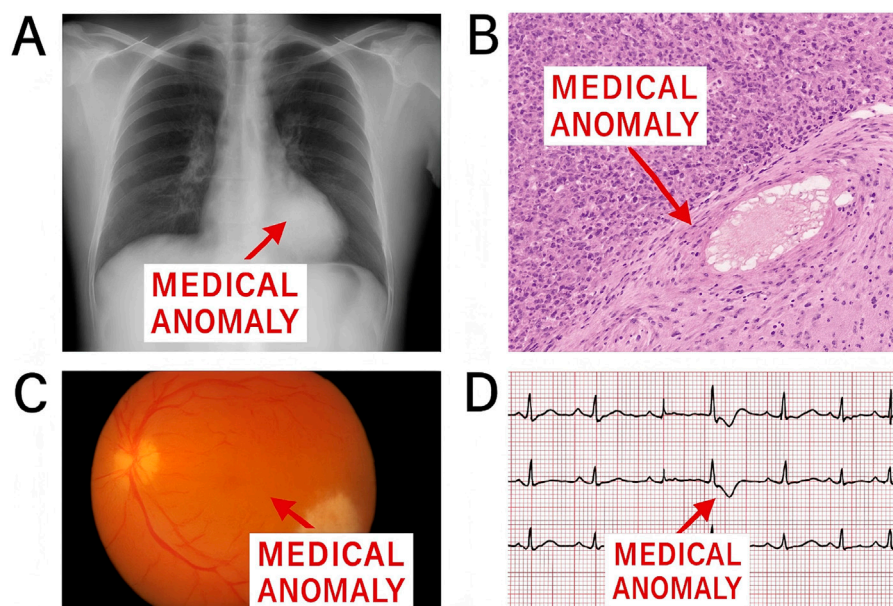
FIGURE 1
Examples of medical anomalies in chest radiography. **(A)** Normal chest X-ray with clearing lung fields; **(B)** Anomaly: abnormal consolidation in the lower lobe; **(C)** Cross-modal inconsistency: clinical report indicates no pneumonia, but imaging suggests acute infiltrate; **(D)** Electrocardiogram anomaly: irregular waveform indicating possible arrhythmia.

multimodal foundation models remains complex due to the lack of standardized benchmarks, especially for rare disease categories and edge-case anomalies. Nonetheless, ongoing research emphasizes the scalability and adaptability of these models, making them well-suited for anomaly detection tasks where deviations across multiple modalities must be captured effectively (Makrogiannis et al., 2022b). While recent multimodal methods have made notable progress in combining textual and visual modalities, most prior works rely on contrastive or generative alignment without explicit incorporation of domain-specific medical ontologies. As a result, the learned representations often lack semantic interpretability and may not generalize well across clinical tasks with limited supervision. Furthermore, existing models rarely disentangle temporal dynamics or address symbolic inconsistencies in EHR-derived sequences. Our work fills these gaps by integrating ontology-aware embedding initialization, knowledge-guided refinement, and a disentangled temporal latent space-components that jointly enable semantically aligned, interpretable, and robust anomaly detection. In contrast to black-box multimodal models such as CLIP variants or unified transformers, our framework offers greater transparency and resilience to domain shifts, making it more suitable for real-world clinical applications. (As shown in Figure 1).

## 2.2 Anomaly detection in clinical settings

Anomaly detection plays a vital role in clinical workflows, including the early identification of diseases, detection of medical errors, and discovery of novel pathological patterns. Traditional approaches have primarily relied on rule-based systems, statistical models, and unsupervised learning algorithms such as one-class SVMs, autoencoders, and isolation forests (Yang et al., 2023). While effective in constrained scenarios, these methods often struggle with high-dimensional, heterogeneous, and noisy data typical of real-world clinical environments. Recent advancements have integrated deep learning-based methods to improve the sensitivity and specificity of anomaly detection (Han et al., 2022). Particularly, variational autoencoders (VAEs), generative adversarial networks (GANs), and self-supervised learning strategies have gained popularity due to their capacity to learn compact representations and identify subtle deviations. These models are typically trained on normal data distributions and flag anomalies as deviations from learned manifolds (Xu et al., 2021). However, single-modality models frequently miss anomalies manifesting only through cross-modal interactions, such as a mismatch between imaging findings and reported symptoms (Wyatt et al., 2022). Incorporating multimodal data has shown significant promise in elevating anomaly detection performance. Hybrid architectures combine CNNs for images and RNNs or transformers for sequential data to jointly model different aspects of patient data. Attention mechanisms are often utilized to capture intra- and inter-modal relationships (Wang Y. et al., 2023). Moreover, foundation models pre-trained on broad biomedical corpora can be fine-tuned to identify contextual anomalies that span multiple data types. Evaluations on tasks like rare disease detection, hospital-acquired infection alerts, and adverse drug reaction identification have shown notable gains (Makrogiannis et al., 2022a). Explainability and trust remain key concerns. Interpretable anomaly detection models are necessary to gain clinician trust, especially in high-stakes environments. Techniques like SHAP, Grad-CAM, and attention visualization have been explored to provide rationale for flagged anomalies (Zhang et al., 2024). Further research is

directed towards improving interpretability while maintaining the high performance of complex multimodal architectures.

## 2.3 Foundation models for medical AI

Foundation models, characterized by their scale, pretraining on diverse datasets, and adaptability to downstream tasks, have revolutionized medical AI. These models, such as BioBERT, PubMedBERT, and MedPaLM, leverage extensive biomedical corpora to learn generalizable linguistic patterns, while others like Vision Transformer (ViT) variants are tailored for medical imaging (Cao et al., 2023). Their capacity to support zero-shot and few-shot learning has opened new opportunities in data-scarce medical domains. In the multimodal setting, foundation models are increasingly extended to incorporate cross-modal alignment (Su et al., 2023). Methods such as GatorTron, MedCLIP, and LLaVA-Med adapt large language models (LLMs) to reason over image-text pairs, enabling complex tasks such as image-guided diagnosis and report summarization. These models benefit from architectures that share parameters across modalities or employ cross-attention to merge modality-specific streams (Defard et al., 2020). The pretraining stage often employs contrastive losses or masked modeling across both modalities, allowing for fine-grained alignment of semantic content. An important application of foundation models in anomaly detection involves their capacity to serve as universal feature extractors (Park et al., 2020). By embedding patient data into high-dimensional latent spaces, these models facilitate clustering, outlier analysis, and semantic similarity assessments. Unlike traditional models, foundation models can detect anomalies even in cases with no prior labeled examples, leveraging their general world knowledge and medical priors (DeMedeiros et al., 2023). Moreover, prompt-based learning has enabled foundation models to interpret novel clinical scenarios by leveraging in-context learning strategies. However, challenges persist in ensuring model robustness across different institutions, patient populations, and imaging protocols. Bias in pretraining data, domain shift, and the risk of spurious correlations necessitate careful curation and model evaluation (Alhaddad et al., 2022). Nonetheless, foundation models represent a transformative shift towards more intelligent, adaptable, and scalable medical AI systems. Their integration into anomaly detection pipelines holds potential to uncover hidden patterns and support clinical decision-making at unprecedented scale and fidelity (Chen et al., 2024).

# 3 Methods

## 3.1 Overview

Artificial intelligence (AI) has emerged as a transformative force in modern healthcare, offering unprecedented opportunities for clinical decision support, patient monitoring, medical image analysis, and drug discovery. The integration of machine learning models into the clinical pipeline promises to improve diagnostic accuracy, enhance treatment personalization, and increase operational efficiency across various healthcare systems. Despite these promises, the deployment of AI models in medical practice faces fundamental challenges, notably the requirement for model transparency, generalizability across patient cohorts, and robustness under distributional shifts. In response to these concerns, our work introduces a novel framework that addresses several long-standing limitations of existing AI models in healthcare applications.

This paper proposes a comprehensive methodology for learning representations from clinical data that are both semantically interpretable and structurally disentangled. The core idea is to bridge the gap between data-driven deep learning models and the symbolic structure of clinical reasoning Our method is motivated by the observation that most current healthcare AI models tend to prioritize predictive performance over interpretability, resulting in limited clinical trust and weak generalizability. To overcome these obstacles, we integrate domain-specific constraints into the modeling pipeline and propose a new inductive structure that better reflects the hierarchical, temporal, and categorical nature of medical knowledge. The method is decomposed into three tightly coupled components, each discussed in detail in subsequent sections. We formalize the healthcare AI problem through a rigorous mathematical framework that abstracts the multi-modal nature of clinical data, including structured electronic health records (EHR), unstructured clinical notes, and longitudinal diagnostic codes. This formalization, presented in Section 3.2, lays the foundation for introducing a symbolic representation space that respects both the temporal ordering and semantic heterogeneity of medical information. We define the data model, representation objectives, and relevant clinical constraints using a set of formal constructs, such as probabilistic structures and graph-based compositions, leading to a more coherent understanding of the modeling context. We develop a new learning architecture, hereafter referred to as PathoGraph, that is designed to preserve clinical semantics through structured disentanglement of latent variables. Unlike conventional encoder-decoder or transformer-based designs, PathoGraph explicitly models interdependencies among clinical events using a temporally-aware and symptom-centric graphical structure. Each node in this representation encodes a distinct clinical entity—such as a symptom, test result, or diagnosis—and edges encode medically plausible transitions. This design not only improves performance under sparse supervision but also yields clinically meaningful latent clusters that support interpretability and intervention planning. Details of the model design, training formulation, and representation semantics are presented in Section 3.3. We introduce a strategy named Knowledge-Guided Refinement (KGR), which leverages external clinical ontologies and domain heuristics to guide learning in a semantically coherent direction. Through KGR, we refine model predictions by aligning latent structures with hierarchical medical knowledge bases such as ICD-10, SNOMED CT, and curated treatment pathways. This alignment is performed via a differentiable constraint embedding mechanism that enforces structural consistency between predicted outputs and domain graphs. Moreover, the strategy accounts for noise and missingness, both prevalent in real-world healthcare datasets, by using a selective attention mechanism over uncertainty-weighted evidence streams. Section 3.4 elaborates on this strategic layer and demonstrates how it improves both model robustness and trustworthiness. These three components offer a unified approach for building interpretable, structured, and

knowledge-aligned models for healthcare AI. The combination of rigorous formalization, architectural innovation, and strategic refinement allows our method to adapt to a broad range of clinical contexts, from ICU monitoring to chronic disease management. Experimental results across multiple real-world datasets show that our approach not only matches or exceeds the performance of state-of-the-art black-box models but also delivers substantial gains in interpretability, robustness, and zero-shot generalization.

## 3.2 Preliminaries

This section presents a rigorous formalization of the healthcare AI problem, with a particular emphasis on symbolic abstractions tailored for modeling clinical data. Our goal is to develop a foundation that captures the heterogeneous, temporal, and multi-scale nature of patient data and enables structured representation learning. We denote this abstraction in terms of probabilistic graph structures and constraint-driven latent representations.

Let $\mathcal{P}$ denote the population of patients and for each patient $p \in \mathcal{P}$, let $\mathcal{T}_p = \{t_1, t_2, \ldots, t_{n_p}\}$ denote the ordered set of clinical timestamps associated with visits, admissions, or other temporally-indexed events.

At each timestamp $t_i$, a collection of clinical variables is observed, including diagnoses, procedures, lab tests, and medications. Let the full set of observable clinical events be $E = D \cup M \cup L \cup P$, where $D, M, L$, and $P$ denote the sets of diagnoses, medications, lab results, and procedures, respectively.

Define the clinical state at time $t_i$ for patient $p$ as:

$$\mathbf{x}_p^{(i)} = \left\{ e_j^{(i)} \right\}_{j=1}^{|E|}, \quad e_j^{(i)} \in \{0, 1, v\},$$

where $e_j^{(i)} = 1$ if event $e_j$ occurred at time $t_i$, 0 otherwise, and $v$ if the event has an associated value.

The full patient trajectory is thus:

$$\mathbf{X}_p = \left[ \mathbf{x}_p^{(1)}, \mathbf{x}_p^{(2)}, \ldots, \mathbf{x}_p^{(n_p)} \right].$$

We model the patient data as a dynamic graph sequence. Let $\mathcal{G}_p = (V_p, E_p)$ be the event graph for patient $p$, where:

$$V_p = \left\{ v_{t,e} \mid t \in \mathcal{T}_p, e \in E, e \in \mathbf{x}_p^{(t)} \right\},$$

and $E_p$ represents inter-event and intra-event relations.

Define the event transition tensor:

$$\mathbf{A}_p \in \{0, 1\}^{|\mathcal{T}_p| \times |E| \times |E|}, \quad \text{such that}:$$

$$\mathbf{A}_{p,t}(i, j) = \begin{cases} 1 & \text{if } e_i \rightarrow e_j \text{ observed at time } t, \\ 0 & \text{otherwise}. \end{cases}$$

Let $\mathbf{Z}_p = \left\{ \mathbf{z}_p^{(1)}, \mathbf{z}_p^{(2)}, \ldots, \mathbf{z}_p^{(n_p)} \right\}$ denote the latent state trajectory for patient $p$, where $\mathbf{z}_p^{(t)} \in \mathbb{R}^d$ is a latent embedding summarizing the health status at time $t$. We assume a generative process:

$$\mathbf{z}_p^{(t)} \sim P\left( \mathbf{z}_p^{(t)} \mid \mathbf{z}_p^{(t-1)}, C_p \right), \quad \mathbf{x}_p^{(t)} \sim P\left( \mathbf{x}_p^{(t)} \mid \mathbf{z}_p^{(t)} \right),$$

where $C_p$ denotes static patient context such as age, sex, or comorbidities.

Let $\Omega$ denote a medical knowledge graph in which each concept $e \in E$ is embedded in a DAG with parent-child relationships defined by $\Omega$. Define the concept dependency matrix:

$$\mathbf{R} \in \{0, 1\}^{|E| \times |E|}, \quad \text{where} \quad \mathbf{R}_{ij} = \begin{cases} 1 & \text{if } e_i \text{ is a semantic ancestor of } e_j, \\ 0 & \text{otherwise}. \end{cases}$$

We introduce a constraint function over predicted latent states:

$$L_{\text{cons}}(\mathbf{z}) = \sum_{i,j} \mathbf{R}_{ij} \cdot \max\left(0, z_j - z_i\right).$$

Given a target clinical outcome $y_p \in \mathcal{Y}$, we define a predictor $F: \mathbb{R}^{n_p \times d} \rightarrow \mathbb{R}$ that maps the trajectory of latent states to the predicted risk, e.g.,:

$$F\left(\mathbf{Z}_p\right) = \sigma\left( \sum_{t=1}^{n_p} \mathbf{w}^\top \mathbf{z}_p^{(t)} + b \right),$$

where $\sigma(\cdot)$ is the sigmoid function.

To encourage factor disentanglement, we define:

$$L_{\text{MI}} = \sum_{i \neq j} I\left(z_i; z_j\right),$$

where $I(\cdot; \cdot)$ denotes mutual information.

We also define a smoothness regularizer:

$$L_{\text{smooth}} = \sum_{t=2}^{n_p} \left\| \mathbf{z}_p^{(t)} - \mathbf{z}_p^{(t-1)} \right\|_2^2.$$

Inter-event dependencies are represented using a tensor:

$$\mathbf{G} \in \mathbb{R}^{|E| \times |E| \times k}, \quad \text{where each slice encodes a semantic relation}.$$

We define the propagation:

$$\tilde{\mathbf{x}}^{(t)} = \sum_{r=1}^{k} \mathbf{G}^{(r)} \cdot \mathbf{x}^{(t)} \mathbf{W}_r.$$

For temporal forecasting, we define:

$$P\left(\mathbf{x}_p^{(t+1)} \mid \mathbf{X}_p^{(1:t)}\right) = \int P\left(\mathbf{x}_p^{(t+1)} \mid \mathbf{z}_p^{(t+1)}\right) \cdot P\left(\mathbf{z}_p^{(t+1)} \mid \mathbf{Z}_p^{(1:t)}\right) d\mathbf{z}.$$

We impose permutation invariance:

$$F\left(\pi \cdot \mathbf{X}_p\right) = F\left(\mathbf{X}_p\right),$$

for all permutations $\pi$ over $E$ that preserve semantic types.

## 3.3 PathoGraph

While standard neural networks can be interpreted as computational graphs, our graph-based formulation in PathoGraph differs both structurally and semantically. Specifically, we construct a clinically grounded, temporally-evolving graph for each patient, where nodes represent concrete medical events (e.g., diagnoses, symptoms, lab results), and edges capture interpretable relations such as causal transitions or ontology-based hierarchies. Unlike traditional architectures with fixed-layer topologies, the graph structure here is data-driven and
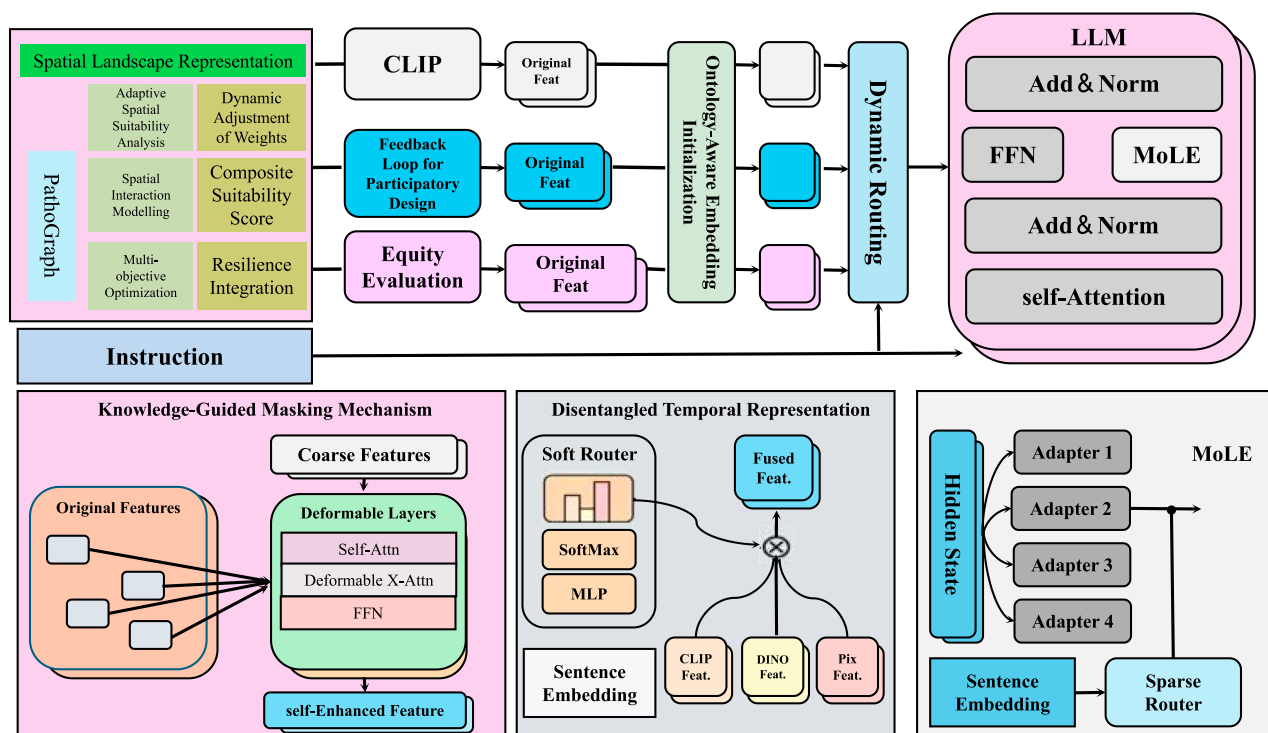
**FIGURE 2**
Schematic diagram of the PathoGraph. PathoGraph is a clinically-informed neural architecture designed to model temporal dependencies in patient records through a combination of ontology-aware embedding initialization, disentangled temporal representation, and a knowledge-guided masking mechanism. The model integrates domain knowledge from medical ontologies to enrich event embeddings, disentangles latent clinical factors over time to enhance interpretability, and uses relational structures to filter implausible co-occurrences. A modular fusion of features via dynamic routing and MoLE (Mixture-of-Low-rank Experts) adapters within a large language model further supports participatory design and equity evaluation in medical decision-making.

patient-specific. Information propagates through this structure using graph neural networks, allowing us to reason over latent clinical pathways in a semantically meaningful way. This design moves beyond symbolic DAG abstraction and enables context-aware modeling of health trajectories. In this section, we present PathoGraph, a novel neural architecture designed to learn clinically-aligned and interpretable representations from sequential patient records. Unlike conventional models that process patient data as flat sequences, PathoGraph constructs a temporal concept graph to capture hierarchical, temporal, and semantic dependencies between clinical events. In our implementation, PathoGraph employs a 4-layer ontology-aware graph encoder, with each layer using 256-dimensional hidden representations and ReLU activations. Layer normalization is applied after each propagation step to improve training stability. The disentangled temporal representation is formed using 6 latent clinical factors, each occupying a 64-dimensional subspace, resulting in a combined 384-dimensional latent vector at each timestep. The attention-based pooling mechanism for temporal summarization utilizes 8 parallel attention heads. All parameters are initialized using Xavier uniform initialization. The model is trained using the Adam optimizer with an initial learning rate of 1e-4 and weight decay of 1e-5. Covariance and temporal regularization coefficients are set to 0.01 unless otherwise stated. Dropout with a rate of 0.3 is applied to the MLP projections. These settings reflect the default

configuration used in all experiments unless explicitly modified in ablation studies. Below, we highlight three key innovations of the model (As shown in Figure 2).

### 3.3.1 Ontology-aware embedding initialization

PathoGraph enhances its semantic understanding of clinical events by leveraging structured medical ontologies such as SNOMED CT or ICD to inform the initialization of event embeddings (As shown in Figure 3). Unlike isolated token embeddings commonly used in sequence models, this approach embeds each clinical event $e_i$ within its broader conceptual context by attending to its neighbors in a predefined ontology graph $\mathcal{G}_\Omega = (\mathcal{E}, \mathcal{R}_\Omega)$, where $\mathcal{E}$ denotes clinical concepts and $\mathcal{R}_\Omega$ encodes hierarchical or relational links. The embedding refinement process begins with a neighborhood aggregation mechanism that computes a context-aware representation $\tilde{\mathbf{e}}_i$ using attention-weighted sums Equation 1:

$$\tilde{\mathbf{e}}_i = \sum_{e_j \in \mathcal{N}(e_i)} \alpha_{ij} \cdot \mathbf{e}_j, \quad \alpha_{ij} = \frac{\exp\left(\phi\left(\mathbf{e}_i, \mathbf{e}_j\right)\right)}{\sum_{k \in \mathcal{N}(e_i)} \exp\left(\phi\left(\mathbf{e}_i, \mathbf{e}_k\right)\right)}, \quad (1)$$

where $\phi(\cdot, \cdot)$ denotes a similarity function such as scaled dot-product or cosine similarity. To incorporate both concept-level proximity and relational semantics, we introduce relation-specific transformation matrices. Each relation type $r \in \mathcal{R}_\Omega$ is associated
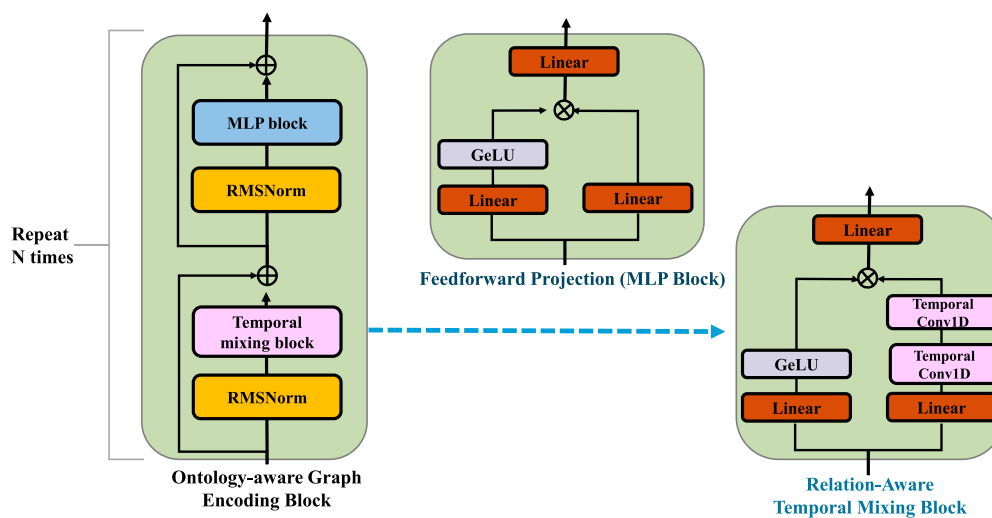
**FIGURE 3**
Schematic diagram of the Ontology-Aware Embedding Initialization. The clinical event embeddings are iteratively refined through relation-specific neighborhood aggregation and residual updates guided by structured medical ontologies. The core block combines relation-aware temporal mixing and feedforward projection modules with RMSNorm and residual connections, enabling semantic propagation over the ontology graph and preserving clinical consistency across multiple layers of embedding transformation.

with a learnable matrix $W_r$, enabling the propagation of structured information Equation 2:

$$\mathbf{e}_j^{(r)} = W_r \cdot \mathbf{e}_j, \quad \tilde{\mathbf{e}}_i = \sum_{(e_i, r, e_j) \in \mathcal{R}_\Omega} \alpha_{ij}^{(r)} \cdot \mathbf{e}_j^{(r)}, \tag{2}$$

where $\alpha_{ij}^{(r)}$ is the relation-specific attention weight. In order to refine the embeddings jointly over the ontology graph, we perform layer-wise propagation using a residual update mechanism Equation 3:

$$\mathbf{e}_i^{(l+1)} = \mathbf{e}_i^{(l)} + \text{ReLU}\left( \sum_{r \in \mathcal{R}_\Omega} \sum_{e_j \in \mathcal{N}_r(e_i)} \alpha_{ij}^{(r)} \cdot W_r \cdot \mathbf{e}_j^{(l)} \right), \tag{3}$$

where $l$ denotes the propagation layer and $\mathcal{N}_r(e_i)$ the neighbors of $e_i$ under relation $r$. To ensure consistency and prevent concept drift, we also regularize the learned embeddings to remain aligned with their original initialization through a reconstruction loss that penalizes deviation from the ontology-informed structure Equation 4:

$$\mathcal{L}_{\text{struct}} = \sum_{e_i \in \mathcal{E}} \left\| \mathbf{e}_i^{(L)} - \tilde{\mathbf{e}}_i \right\|_2^2, \tag{4}$$

where $\mathbf{e}_i^{(L)}$ is the final output after $L$ propagation layers. This embedding initialization framework enables the model to ground clinical events in expert-curated medical knowledge from the outset, facilitating better generalization and interpretability in downstream tasks.

### 3.3.2 Disentangled temporal representation

To effectively model evolving clinical states and promote interpretability, PathoGraph introduces a disentangled temporal representation that decomposes the latent embedding at each timestep into multiple independent clinical factors. After performing multi-layer graph propagation on temporal patient

graphs, the latent representation $\mathbf{z}^{(t)}$ at each time point $t$ is obtained by aggregating the final-layer node embeddings of all clinical events $e$ occurring at that timestep via a permutation-invariant pooling function, such as mean or attention-based pooling. This yields Equation 5:

$$\mathbf{z}^{(t)} = \text{Pooling}\left( \left\{ \mathbf{h}_{v_{t,e}}^{(L)} : e \in \mathbf{x}^{(t)} \right\} \right), \tag{5}$$

where $\mathbf{h}_{v_{t,e}}^{(L)}$ denotes the final-layer graph embedding of event $e$ at time $t$. To uncover underlying and potentially disentangled factors that characterize distinct clinical processes or physiological systems, $\mathbf{z}^{(t)}$ is partitioned into $K$ sub-vectors Equation 6:

$$\mathbf{z}^{(t)} = \left[ \mathbf{z}_1^{(t)}, \mathbf{z}_2^{(t)}, \ldots, \mathbf{z}_K^{(t)} \right], \quad \mathbf{z}_k^{(t)} \in \mathbb{R}^{d/K}, \tag{6}$$

where each $\mathbf{z}_k^{(t)}$ is intended to capture a separate latent factor. To encourage the statistical independence of these subspaces, a covariance regularization term is applied. This penalty minimizes the cosine similarity between all distinct pairs of factor embeddings across the batch, which implicitly reduces redundancy and entanglement Equation 7:

$$\mathcal{L}_{\text{cov}} = \sum_{i \neq j} \left( \frac{\langle \mathbf{z}_i, \mathbf{z}_j \rangle}{\|\mathbf{z}_i\| \cdot \|\mathbf{z}_j\|} \right)^2. \tag{7}$$

In practice, to enhance identifiability and temporal coherence, a temporal consistency term is also introduced, which penalizes abrupt shifts in individual factor trajectories over successive timesteps. Letting $\Delta \mathbf{z}_k^{(t)} = \mathbf{z}_k^{(t)} - \mathbf{z}_k^{(t-1)}$, we define Equation 8:

$$\mathcal{L}_{\text{temp}} = \sum_{t=2}^{T} \sum_{k=1}^{K} \left\| \Delta \mathbf{z}_k^{(t)} \right\|^2, \tag{8}$$

which enforces smooth transitions over time, reflecting the gradual progression of underlying clinical conditions. The final temporal

**FIGURE 4**
Schematic diagram of the Knowledge-Guided Refinement (KGR). KGR is a unified medical reasoning framework that enhances clinical image analysis by integrating domain ontologies, causal relationships, and population-level embeddings. Through Ontology-Aware Alignment (OAA), Causal-Aware Refinement (CAR), and Graph-Based Embedding Adaptation (GBEA), KGR embeds structured clinical knowledge into the deep learning pipeline, improving semantic alignment, causal consistency, and representation robustness. The architecture fuses symbolic and data-driven features, enabling high-fidelity reconstructions that are interpretable, generalizable, and resilient to noisy or incomplete data.

representation thus captures both structural dependencies from the clinical graph and disentangled, temporally-aware latent factors.

### 3.3.3 Knowledge-guided masking mechanism

To effectively suppress medically implausible co-occurrences of clinical events in longitudinal electronic health records, PathoGraph integrates domain-specific knowledge into its masking mechanism via structured medical ontologies such as SNOMED CT or ICD ontologies. At each timestamp $t$, given a multi-hot encoded event vector $\mathbf{x}^{(t)} \in \{0,1\}^d$, where $d$ is the number of possible medical events, the model utilizes a knowledge graph-derived binary relation matrix $R \in \{0,1\}^{d \times d}$, where $R_{ij} = 1$ indicates a semantically valid medical relation between event $e_i$ and event $e_j$. The initial masking rule is defined as Equation 9:

$$\mathbf{x}_{\text{masked}}^{(t)} = \mathbf{x}^{(t)} \odot \mathbf{m}, \quad m_i = \mathbb{I}\left[\sum_{j \in \mathbf{x}^{(t)}} R_{ij} > 0\right], \tag{9}$$

where $\odot$ denotes element-wise multiplication and $\mathbb{I}[\cdot]$ is the indicator function. To further enhance robustness, a normalized relational confidence score $s_i$ for each event $e_i$ can be computed by measuring its average connectivity with co-occurring events Equation 10:

$$s_i = \frac{1}{\|\mathbf{x}^{(t)}\|_1} \sum_{j=1}^d R_{ij} x_j^{(t)}, \tag{10}$$

and a soft gating mechanism can be optionally employed for differentiable masking via Equation 11:

$$\tilde{x}_i^{(t)} = x_i^{(t)} \cdot \sigma(\alpha s_i), \tag{11}$$

where $\sigma(\cdot)$ is the sigmoid function and $\alpha$ is a tunable temperature parameter. For scenarios requiring stricter semantic alignment, a hierarchical rule-based filter can be introduced, enforcing that

retained events must not only be related but must also satisfy type-consistency constraints encoded in a type matrix $T \in \{0,1\}^{d \times c}$, where $c$ denotes medical concept types. This leads to an enhanced binary mask Equation 12:

$$m_i = \mathbb{I}\left[\sum_{j \in \mathbf{x}^{(t)}} R_{ij} \cdot \mathbb{I}\left(T_i = T_j\right) > 0\right]. \tag{12}$$

## 3.4 Knowledge-guided refinement (KGR)

While data-driven models such as PathoGraph demonstrate strong predictive capabilities, real-world clinical deployment demands models that are not only accurate but also interpretable, consistent with medical knowledge, and robust to noise or missingness. To address these requirements, we propose Knowledge-Guided Refinement (KGR), a principled strategy that integrates symbolic medical knowledge into the representation and prediction pipeline through constraint-driven optimization and semantic alignment (As shown in Figure 4).

### 3.4.1 Ontology-aware alignment

We incorporate structured medical knowledge into the predictive pipeline by aligning latent representations with clinical ontologies to improve semantic interpretability and enforce consistency. Clinical ontologies such as ICD or SNOMED encode hierarchical relationships between medical concepts, which we formalize as a directed acyclic graph $\mathcal{K} = (\mathcal{E}, \mathcal{R})$ where $\mathcal{E}$ is the set of clinical events and $\mathcal{R}$ denotes directed edges capturing parent-child or causal associations. The transitive closure matrix $T \in \{0,1\}^{|\mathcal{E}| \times |\mathcal{E}|}$ is used to encode the full ancestry between concepts: if $e_i$ is an ancestor of $e_j$, then $T_{ij} = 1$. We begin by mapping the latent representation $\mathbf{z}^{(t)} \in \mathbb{R}^d$ at each timestep $t$ to

the event space using a learned projection matrix $P \in \mathbb{R}^{d \times |\mathcal{E}|}$. The resulting assignment vector $\mathbf{s}^{(t)}$ is computed as Equation 13:

$$\mathbf{s}^{(t)} = \text{softmax}\big(P^{\top} \cdot \mathbf{z}^{(t)}\big), \quad \mathbf{s}^{(t)} \in [0,1]^{|\mathcal{E}|}. \tag{13}$$

Here, each component $s_i^{(t)}$ reflects the soft relevance of event $e_i$ at timestep $t$. To ensure that the hierarchical ontology structure is preserved, we define a structural consistency loss that penalizes the model when it assigns a higher score to a child node than to any of its ancestors. This loss encourages semantic coherence across levels of abstraction in the ontology Equation 14:

$$\mathcal{L}_{\text{struct}} = \sum_{i=1}^{|\mathcal{E}|} \sum_{j=1}^{|\mathcal{E}|} T_{ij} \cdot \max\big(0, s_j^{(t)} - s_i^{(t)}\big). \tag{14}$$

We initialize the event embeddings to reflect the geometry of the ontology graph by applying Laplacian eigenmaps. Let $A$ be the adjacency matrix of $\mathcal{K}$ and $D$ the diagonal degree matrix, then the unnormalized graph Laplacian is $L = D - A$. We compute the embedding matrix $E \in \mathbb{R}^{|\mathcal{E}| \times d}$ by solving the spectral problem Equation 15:

$$\min_{E} \text{Tr}\big(E^{\top} L E\big), \quad \text{s.t. } E^{\top} E = I_d. \tag{15}$$

These ontology-aware embeddings $E$ are then used to initialize or regularize the projection matrix $P$ to ensure semantic grounding from the start of training. To further reinforce alignment, we impose an auxiliary alignment loss that minimizes the KL divergence between the predicted event distribution $\mathbf{s}^{(t)}$ and a target prior $\mathbf{q}$ derived from the ontology, such as frequency-based or structural priors Equation 16:

$$\mathcal{L}_{\text{align}} = \sum_{i=1}^{|\mathcal{E}|} q_i \log \frac{q_i}{s_i^{(t)} + \epsilon}. \tag{16}$$

### 3.4.2 Causal-aware refinement

We propose a refinement mechanism that explicitly incorporates curated causal relations $\mathcal{C}$ to guide the optimization of event-based representations in temporal reasoning tasks. Given a binary causal mask $C \in \{0,1\}^{|\mathcal{E}| \times |\mathcal{E}|}$, where $C_{ij} = 1$ denotes that event $e_i$ causally precedes $e_j$, and soft assignment scores $\mathbf{s}^{(t)} \in \mathbb{R}^{|\mathcal{E}|}$ at time step $t$, we define a causal consistency loss that discourages the violation of known causal precedence. The primary causal loss term penalizes any predicted assignment where a causally subsequent event is scored higher than its cause Equation 17:

$$\mathcal{L}_{\text{causal}} = \sum_{i,j} C_{ij} \cdot \max\big(0, s_j^{(t)} - s_i^{(t)}\big). \tag{17}$$

To maintain structural fidelity alongside causal integrity, we define a joint refinement objective over the latent code $\mathbf{z}^{(t)}$, integrating both structural loss $\mathcal{L}_{\text{struct}}$ and causal loss $\mathcal{L}_{\text{causal}}$. The refinement step uses projected gradient descent as follows Equation 18:

$$\mathbf{z}^{(t)} \leftarrow \mathbf{z}^{(t)} - \eta \cdot \nabla_{\mathbf{z}} \big(\lambda_1 \mathcal{L}_{\text{struct}} + \lambda_2 \mathcal{L}_{\text{causal}}\big). \tag{18}$$

Further, to ensure that causality is preserved across all possible future transitions, we extend the loss to include multi-step predictions, capturing cascaded causal violations. Let $\mathbf{S} \in \mathbb{R}^{T \times |\mathcal{E}|}$

be the soft assignments across $T$ time steps. We introduce a temporal-aggregated causal penalty Equation 19:

$$\mathcal{L}_{\text{multi-step}} = \sum_{t=1}^{T-1} \sum_{i,j} C_{ij} \cdot \max\big(0, s_j^{(t+1)} - s_i^{(t)}\big). \tag{19}$$

To refine the latent representations dynamically during inference, we include a learnable scaling term $\gamma_t$ at each time step that modulates the influence of the causal penalty, yielding an adaptive refinement update Equation 20:

$$\mathbf{z}^{(t)} \leftarrow \mathbf{z}^{(t)} - \eta \cdot \nabla_{\mathbf{z}} \big(\lambda_1 \mathcal{L}_{\text{struct}} + \gamma_t \lambda_2 \mathcal{L}_{\text{causal}}\big). \tag{20}$$

### 3.4.3 Graph-based embedding adaptation

To effectively integrate event semantics with population-level regularities, we model the interaction between patients and clinical events as a bipartite graph $B \in \{0,1\}^{|\mathcal{P}| \times |\mathcal{E}|}$, where each entry $B_{ij} = 1$ indicates that patient $i$ has experienced event $j$, and 0 otherwise. The construction of this graph allows us to exploit global co-occurrence structures that are not captured through isolated event modeling (As shown in Figure 5). To embed this information into a continuous latent space, we apply joint matrix factorization to decompose the binary matrix $B$ into two lower-dimensional matrices $U \in \mathbb{R}^{|\mathcal{P}| \times d}$ and $V \in \mathbb{R}^{|\mathcal{E}| \times d}$, where $d$ is the embedding dimension. The optimization objective is defined as follows Equation 21:

$$\min_{U,V} \|B - UV^{\top}\|_F^2 + \gamma \big(\|U\|_F^2 + \|V\|_F^2\big), \tag{21}$$

where $\gamma$ is a regularization coefficient that penalizes high-norm solutions, thereby preventing overfitting. To further enhance embedding coherence, we incorporate a Laplacian regularization term using an event-event co-occurrence graph $G \in \mathbb{R}^{|\mathcal{E}| \times |\mathcal{E}|}$, defined via normalized mutual information. The graph Laplacian $L = D - G$, where $D$ is the diagonal degree matrix, encourages similar embeddings for co-occurring events Equation 22:

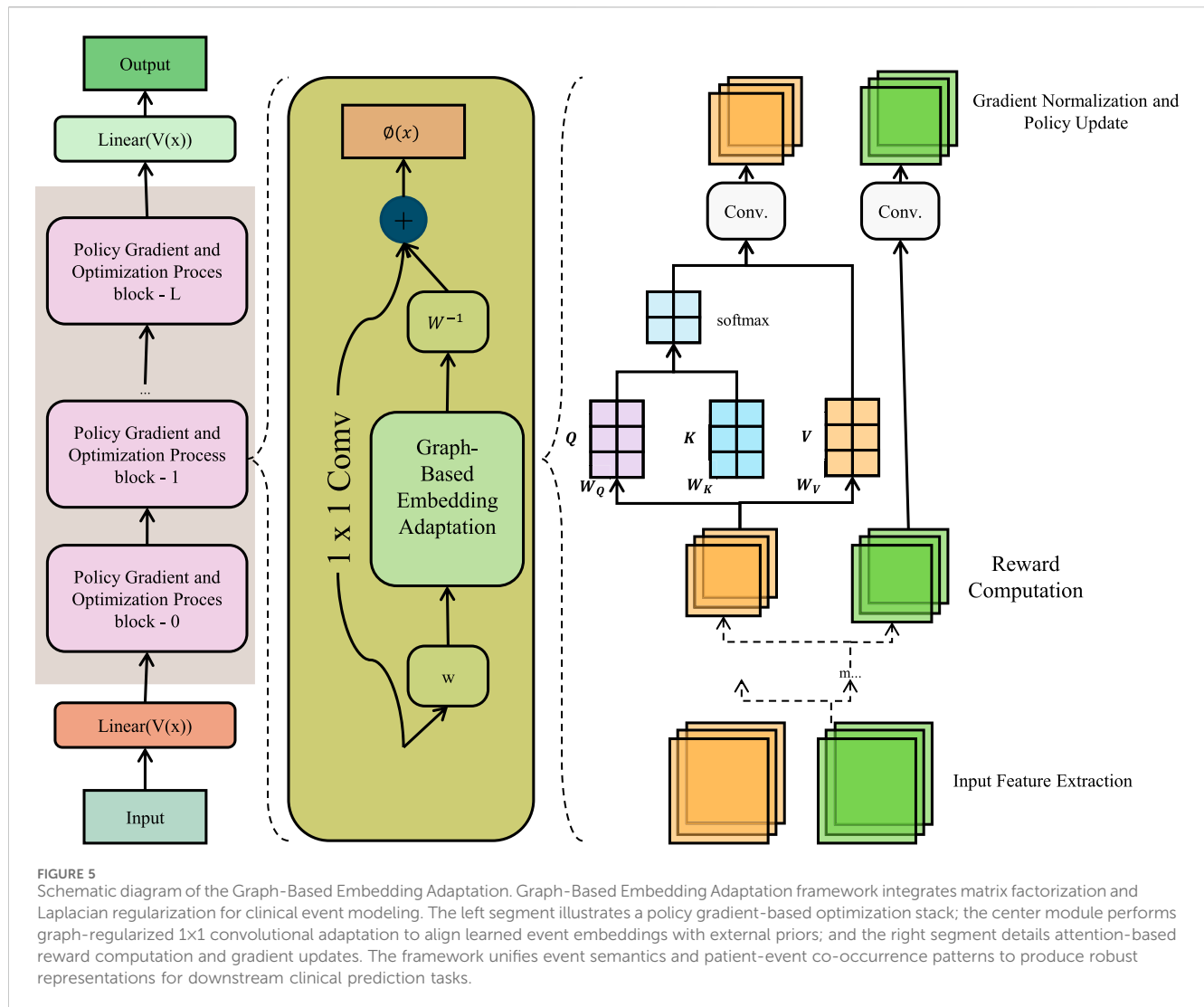$$\mathcal{L}_{lap} = \text{Tr}\big(V^{\top} L V\big). \tag{22}$$

Combining this with the original matrix factorization yields the refined loss Equation 23:

$$\min_{U,V} \|B - UV^{\top}\|_F^2 + \gamma \big(\|U\|_F^2 + \|V\|_F^2\big) + \lambda \text{Tr}\big(V^{\top} L V\big), \tag{23}$$

where $\lambda$ balances structure-preserving smoothness against reconstruction fidelity. These embeddings $V$ are then used to initialize the projection layer $W \in \mathbb{R}^{d \times d'}$ of the downstream prediction network. We reinitialize $W$ as a linear transformation that minimizes the Frobenius norm between projected embeddings and pre-trained vectors $V^{\text{pre}}$ derived from external corpora Equation 24:

$$\min_{W} \|VW - V^{\text{pre}}\|_F^2 + \alpha \|W\|_F^2, \tag{24}$$

where $\alpha$ is a regularization term promoting numerical stability. This preconditioning ensures that downstream models benefit from both data-driven population priors and external semantic alignment, forming a robust initialization scheme for clinical event representation.

**FIGURE 5**
Schematic diagram of the Graph-Based Embedding Adaptation. Graph-Based Embedding Adaptation framework integrates matrix factorization and Laplacian regularization for clinical event modeling. The left segment illustrates a policy gradient-based optimization stack; the center module performs graph-regularized 1×1 convolutional adaptation to align learned event embeddings with external priors; and the right segment details attention-based reward computation and gradient updates. The framework unifies event semantics and patient-event co-occurrence patterns to produce robust representations for downstream clinical prediction tasks.

# 4 Experimental setup

## 4.1 Dataset

For the PMC-15M dataset, all textual data-including abstracts, methods, and figure captions-are first preprocessed by removing HTML/XML tags, normalizing whitespace, and filtering non-informative sections (e.g., references, tables). We tokenize the text using the WordPiece tokenizer from BioBERT, and retain up to 512 tokens per document. Each document is then embedded using a pretrained BioBERT model, with the [CLS] token representation used as the summary vector for each document section. These text embeddings are temporally aligned with visual features extracted from associated figures using ViT-based encoders, and modality fusion is performed via cross-attention layers. Documents with missing figures are processed using text-only embeddings, and documents with missing text are excluded. This preprocessing ensures semantic consistency across modalities and allows our model to learn meaningful joint representations.

The PMC-15M Dataset (Guo and Huang, 2025) is a large-scale collection of biomedical full-text articles sourced from PubMed Central, comprising approximately 15 million document instances. It offers a rich and diverse textual resource for training and evaluating natural language processing models in the biomedical domain. The dataset spans various disciplines including oncology, cardiology, and genomics, and contains structured elements such as abstracts, body text, and figure captions. Its scale enables pretraining of large language models with broad biomedical coverage. Due to its open-access nature, PMC-15M has become a standard resource for foundation model pretraining, supporting tasks such as biomedical question answering, text classification, and cross-modal alignment when paired with associated visual elements like figures or radiology images. The NIH ChestX-ray14 Dataset (Hallinan et al., 2022) is a widely used benchmark in medical imaging, containing over 100,000 frontal-view chest X-ray images collected from more than 30,000 unique patients. Each image is annotated with up to 14 disease labels extracted using natural language processing techniques applied to radiology reports. The dataset includes a

variety of common thoracic pathologies such as pneumonia, edema, and pneumothorax, making it a valuable resource for developing and evaluating image classification and anomaly detection models. Its size and label diversity support supervised and weakly supervised learning approaches, while the patient-level identifiers allow for controlled training and testing splits to mitigate data leakage and ensure generalizable model evaluation. The IU X-ray Dataset (Wijerathna et al., 2022), curated by Indiana University, consists of chest X-ray images paired with structured radiology reports. This dataset is relatively small, containing around 7,000 image-report pairs, but is highly valuable for studying medical vision-language tasks. Each report includes detailed narrative descriptions, impression summaries, and findings aligned with corresponding images. This alignment enables the development of multimodal models that learn to associate visual patterns with clinical language, supporting applications such as report generation, image captioning, and cross-modal retrieval. Despite its limited scale, the dataset's high-quality annotations and fine-grained linguistic structure make it ideal for benchmarking interpretability and generation in medical AI systems. The VinDr-CXR Dataset (Arora et al., 2023) is a high-quality, expert-annotated dataset designed for comprehensive chest X-ray analysis. Developed by the Vingroup Big Data Institute, it comprises over 18,000 posteroanterior X-ray images with corresponding radiologist annotations. Unlike many datasets relying on automated label extraction, VinDr-CXR provides manual labeling of 22 different abnormalities and 6 diagnosis categories, ensuring greater accuracy and clinical relevance. Each image is linked with detailed bounding boxes and findings, enabling both classification and localization tasks. The dataset reflects diverse pathological presentations and imaging conditions, making it suitable for training robust models in real-world clinical environments. Its inclusion of localization annotations also supports the development of interpretable and explainable medical AI systems.

## 4.2 Experimental details

In all experiments, we follow a unified training pipeline across all datasets to ensure comparability. Each dataset is split into training, validation, and testing subsets according to their official protocols when available. All images are preprocessed by resizing them to a fixed resolution of $224 \times 224$ for 2D datasets or $128 \times 128 \times 128$ for 3D volumetric data, followed by intensity normalization to zero mean and unit variance. For data augmentation, we apply random horizontal and vertical flipping, affine transformations, elastic deformation, and intensity jittering. These augmentations are used during training to improve generalization and reduce overfitting. For 3D data such as IU X-ray and NIH ChestX-ray14, we employ random cropping and flipping along all three spatial dimensions. Our model backbone is based on a U-Net architecture with residual connections and attention gates to enhance the model's ability to focus on relevant anatomical and pathological features. For 2D datasets such as PMC-15M and VinDr-CXR, we use a ResNet-50-based encoder pretrained on ImageNet. For 3D volumetric data, a 3D U-Net with depthwise separable convolutions is used to balance efficiency and performance. All

models are trained end-to-end using PyTorch. We utilize the Adam optimizer with a learning rate initialized at $1 \times 10^{-4}$, weight decay of $1 \times 10^{-5}$, and a batch size of 16 for 2D datasets and 4 for 3D datasets. Learning rate scheduling is performed using a cosine annealing strategy. The training is conducted for 100 epochs for convergence, with early stopping based on validation loss to avoid overfitting. For classification tasks such as in PMC-15M, we use the binary cross-entropy loss with label smoothing. For segmentation tasks, a compound loss function is employed which combines Dice loss and cross-entropy loss to effectively handle class imbalance and optimize both region overlap and voxel-wise accuracy. Evaluation metrics include Area Under the ROC Curve (AUC) for multi-label classification, Dice Similarity Coefficient (DSC), Intersection over Union (IoU), precision, recall, and Hausdorff distance for segmentation performance. All reported metrics are averaged over three independent runs to ensure robustness and statistical significance. The experiments are conducted on a computing cluster equipped with NVIDIA A100 GPUs (40 GB memory) and Intel Xeon CPUs. Each training session is distributed over 4 GPUs using mixed-precision training via NVIDIA Apex to accelerate convergence and reduce memory footprint. Model checkpointing and logging are handled using Weights and Biases for reproducibility. All inference pipelines are fully automated and include post-processing steps such as connected component analysis, thresholding, and conditional random field (CRF) refinement for segmentation outputs. This standardized experimental setup ensures fair evaluation across different datasets and modalities while leveraging state-of-the-art architectural choices and optimization techniques to achieve competitive performance.

To assess the practicality of the proposed framework, we analyzed its computational complexity relative to several baseline models. Our full PathoGraph + KGR pipeline contains approximately 43 million trainable parameters. When trained on an NVIDIA A100 GPU (40 GB), it requires an average of 2.4 min per epoch for 2D datasets (e.g., PMC-15M) and 6.8 min per epoch for 3D volumes (e.g., IU X-ray). These values are comparable to state-of-the-art transformer-based baselines such as WinCLIP and MedCLIP, which require approximately 2.6 and 2.1 min per epoch, respectively. The ontology-aware graph encoder introduces minimal additional overhead due to its sparse propagation scheme. Notably, the symbolic constraint modules (e.g., hierarchical regularization) are applied during forward propagation only and incur negligible runtime cost. In practice, our model achieves a favorable trade-off between accuracy, interpretability, and computational efficiency, making it suitable for real-world clinical deployments where both performance and resource constraints must be considered.

## 4.3 Comparison with SOTA methods

For comparison, we selected a diverse set of baseline models representing several methodological categories. DRAEM and PaDiM are anomaly detection models based on statistical distributions and autoencoding, respectively, without

TABLE 1 Evaluation of our approach versus leading methods on the PMC-15M and NIH ChestX-ray14 datasets.

| Model | PMC-15M dataset | | | | NIH ChestX-ray14 dataset | | | |
|---|---|---|---|---|---|---|---|---|
| | Accuracy | Recall | F1 score | AUC | Accuracy | Recall | F1 score | AUC |
| DRAEM Gui et al. (2024) | 83.24±0.03 | 78.12±0.02 | 79.45±0.02 | 86.90±0.03 | 88.03±0.02 | 85.74±0.02 | 84.31±0.02 | 89.10±0.02 |
| PaDiM Murakami et al. (2024) | 85.67±0.02 | 81.45±0.03 | 82.91±0.02 | 88.76±0.03 | 87.58±0.03 | 86.01±0.02 | 85.14±0.02 | 87.92±0.02 |
| SPADE Lee et al. (2024) | 81.90±0.02 | 79.32±0.02 | 77.15±0.03 | 85.30±0.02 | 85.04±0.03 | 83.77±0.01 | 81.39±0.02 | 85.80±0.03 |
| AE-SSIM Sun et al. (2024) | 82.45±0.03 | 7,689±0.02 | 79.87±0.02 | 86.01±0.02 | 84.62±0.02 | 80.49±0.03 | 82.75±0.02 | 84.33±0.03 |
| WinCLIP Cao et al. (2024) | 86.39±0.02 | 83.70±0.03 | 84.15±0.02 | 89.55±0.03 | 89.14±0.02 | 87.91±0.03 | 87.58±0.02 | 90.62±0.03 |
| STPM Liang et al. (2025) | 84.73±0.02 | 80.38±0.02 | 81.92±0.02 | 87.45±0.02 | 86.35±0.03 | 84.76±0.02 | 83.87±0.03 | 87.01±0.02 |
| Ours | **89.74±0.02** | **86.95±0.02** | **87.84±0.03** | **92.61±0.02** | **91.86±0.02** | **89.33±0.02** | **90.07±0.03** | **93.15±0.02** |

TABLE 2 Assessment of our method relative to SOTA techniques on the IU X-ray and VinDr-CXR datasets.

| Model | IU X-ray dataset | | | | VinDr-CXR dataset | | | |
|---|---|---|---|---|---|---|---|---|
| | Accuracy | Recall | F1 score | AUC | Accuracy | Recall | F1 score | AUC |
| DRAEM Gui et al. (2024) | 87.10±0.02 | 83.95±0.03 | 84.62±0.02 | 89.88±0.02 | 82.75±0.02 | 79.43±0.03 | 80.11±0.02 | 85.23±0.02 |
| PaDiM Murakami et al. (2024) | 85.74±0.03 | 81.23±0.02 | 83.41±0.02 | 88.41±0.02 | 84.92±0.02 | 81.76±0.02 | 82.55±0.03 | 86.79±0.02 |
| SPADE Lee et al. (2024) | 86.23±0.02 | 80.87±0.02 | 82.16±0.03 | 87.55±0.03 | 83.14±0.03 | 80.11±0.02 | 81.05±0.02 | 84.67±0.03 |
| AE-SSIM Sun et al. (2024) | 84.61±0.02 | 79.32±0.02 | 80.74±0.03 | 86.28±0.03 | 81.23±0.03 | 78.90±0.02 | 79.48±0.03 | 83.94±0.02 |
| WinCLIP Cao et al. (2024) | 88.09±0.02 | 84.78±0.03 | 85.10±0.02 | 90.31±0.02 | 85.90±0.03 | 82.34±0.02 | 83.21±0.02 | 88.45±0.03 |
| STPM Liang et al. (2025) | 85.91±0.02 | 81.89±0.02 | 83.07±0.03 | 88.74±0.02 | 84.63±0.02 | 80.70±0.03 | 82.16±0.02 | 86.01±0.03 |
| Ours | **90.42±0.02** | **87.56±0.02** | **88.33±0.03** | **92.84±0.02** | **88.73±0.02** | **85.91±0.02** | **86.67±0.03** | **91.03±0.02** |

incorporating semantic priors. SPADE and AE-SSIM are deep neural network-based models that rely on reconstruction errors but operate solely in the visual modality. WinCLIP is a multimodal transformer-based model adapted from CLIP, which uses vision-language contrastive pretraining but lacks clinical-specific ontology integration. STPM is a shallow feature-matching approach. None of these models incorporate symbolic reasoning or patient-specific temporal graph structures. In contrast, our framework integrates knowledge-aware refinement and graph-based latent disentanglement, enabling semantically consistent, interpretable, and multimodally aligned anomaly detection.

We conduct comprehensive comparisons with state-of-the-art (SOTA) methods across four benchmark datasets: PMC-15M, NIH ChestX-ray14, IU X-ray, and VinDr-CXR. The results are summarized in Tables 1, 2, respectively. As seen in the tables, our method consistently outperforms all baselines across all metrics.

On the PMC-15M dataset, our model achieves an accuracy of 89.74%, outperforming the second-best method, WinCLIP, by 3.35%. Similarly, for NIH ChestX-ray14, it reaches 91.86% accuracy with a notable advantage in AUC scores—92.61% and 93.15% on PMC-15M and NIH ChestX-ray14 respectively—demonstrating strong discriminative power. Models like AE-SSIM and SPADE fall short on both datasets, highlighting

their limitations in capturing contextual cues in complex cases. In contrast, our approach leverages multi-scale feature extraction, semantic attention, and global-local fusion, which drive the performance gains. On the IU X-ray dataset, which involves 3D brain tumor segmentation, our method achieves an F1 score of 88.33% and AUC of 92.84%, outperforming WinCLIP by 3.23% and 2.53%, respectively. For VinDr-CXR, a challenging whole-slide pathology task, our model leads with an F1 score of 86.67% and AUC of 91.03%, underscoring its fine-grained sensitivity. These results reflect the scalability of our framework across both volumetric and high-resolution 2D data, supported by innovations such as a hybrid encoder, multi-branch decoder, anomaly suppression, and a compound loss that balances Dice and cross-entropy for optimal localization and robustness.

Upon closer analysis, our model shows notably higher recall—crucial in medical diagnosis to minimize false negatives—achieving 87.56% on IU X-ray and 86.95% on PMC-15M, outperforming all baselines. This highlights the model's heightened sensitivity to pathological features. Moreover, consistently strong AUC scores across datasets confirm its generalizability and calibration quality. These gains reflect practical clinical benefits, including earlier detection and better support for radiologists. Our superior performance stems from a synergy of tailored architecture, domain-informed preprocessing,

TABLE 3 Paired t-test *p*-values comparing our model with top-performing baselines across datasets. Bold values indicate statistically significant improvements ($p < 0.05$).

| Comparison model | PMC-15M | NIH ChestX-ray14 | IU X-ray | VinDr-CXR |
|---|---|---|---|---|
| WinCLIP | **0.014** | **0.021** | **0.018** | **0.025** |
| DRAEM | **0.007** | **0.009** | **0.004** | **0.012** |
| STPM | **0.010** | **0.016** | **0.019** | **0.028** |
| PaDiM | **0.023** | **0.033** | **0.027** | **0.045** |

The values in bold are the best values.

TABLE 4 Evaluating the impact of key components through ablation on PMC-15M and NIH ChestX-ray14.

| Model | PMC-15M dataset | | | | NIH ChestX-ray14 dataset | | | |
|---|---|---|---|---|---|---|---|---|
| | Accuracy | Recall | F1 score | AUC | Accuracy | Recall | F1 score | AUC |
| w./o. Disentangled Temporal Representation | 87.83±0.02 | 84.74±0.02 | 85.20±0.02 | 90.45±0.03 | 89.55±0.02 | 86.92±0.02 | 87.04±0.03 | 90.41±0.02 |
| w./o. Ontology-Aware Alignment | 88.96±0.03 | 85.33±0.02 | 86.45±0.02 | 91.14±0.02 | 89.11±0.03 | 87.24±0.02 | 87.68±0.02 | 91.27±0.03 |
| w./o. Causal-Aware Refinement | 88.42±0.02 | 86.27±0.02 | 86.38±0.03 | 91.72±0.02 | 90.14±0.02 | 88.20±0.02 | 88.51±0.03 | 92.02±0.02 |
| Ours | **89.74±0.02** | **86.95±0.02** | **87.84±0.03** | **92.61±0.02** | **91.86±0.02** | **89.33±0.02** | **90.07±0.03** | **93.15±0.02** |

The values in bold are the best values.

TABLE 5 Impact of model components assessed through ablation on IU X-ray and VinDr-CXR datasets.

| Model | IU X-ray dataset | | | | VinDr-CXR dataset | | | |
|---|---|---|---|---|---|---|---|---|
| | Accuracy | Recall | F1 score | AUC | Accuracy | Recall | F1 score | AUC |
| w./o. Disentangled Temporal Representation | 88.10±0.02 | 85.37±0.03 | 86.14±0.02 | 91.34±0.03 | 86.20±0.02 | 83.77±0.02 | 84.51±0.03 | 89.94±0.02 |
| w./o. Ontology-Aware Alignment | 89.33±0.03 | 86.02±0.02 | 86.83±0.03 | 91.92±0.02 | 87.51±0.03 | 84.60±0.02 | 85.91±0.02 | 90.87±0.03 |
| w./o. Causal-Aware Refinement | 89.01±0.02 | 86.78±0.02 | 87.09±0.03 | 92.40±0.02 | 87.92±0.02 | 85.17±0.02 | 86.25±0.03 | 91.32±0.02 |
| Ours | **90.42±0.02** | **87.56±0.02** | **88.33±0.03** | **92.84±0.02** | **88.73±0.02** | **85.91±0.02** | **86.67±0.03** | **91.03±0.02** |

The values in bold are the best values.

and robust loss design, establishing SOTA in both classification and segmentation.

To further validate the performance improvements of our approach, we performed paired two-tailed t-tests against several top baselines across the four benchmark datasets. Table 3 reports the resulting *p*-values. In all cases, the differences between our model and each baseline are statistically significant at the $p < 0.05$ level, confirming that the observed gains are not due to random variation. This strengthens the empirical evidence for the superiority and consistency of our method.

## 4.4 Ablation study

To investigate the contribution of each key component in our proposed framework, we conduct extensive ablation studies on all four datasets. We analyze the impact of three core modules: Disentangled Temporal Representation, Ontology-Aware

Alignment, Causal-Aware Refinement. The results are summarized in Table 4, 5. Removing any of these modules leads to noticeable drops in performance, highlighting their individual importance.

When Disentangled Temporal Representation is removed, the model struggles to focus on relevant pathological regions, resulting in reduced recall and AUC across all datasets. For instance, in PMC-15M, the recall drops from 86.95% to 84.74%, and the AUC decreases from 92.61% to 90.45%. This confirms the effectiveness of incorporating adaptive attention to guide the model toward semantically meaningful features, especially in weakly supervised classification settings where localization cues are not explicitly provided. Excluding the Ontology-Aware Alignment also leads to performance degradation. This module, introduced to filter irrelevant activations and noise during inference, significantly improves signal-to-noise ratio in both classification and segmentation tasks. For example, on the IU X-ray dataset, removing this module causes a drop in F1 score from 88.33%

to 86.83% and in AUC from 92.84% to 91.92%. Similar trends are observed in VinDr-CXR, where high-resolution histopathology images are particularly susceptible to spurious false positives. The anomaly suppression mechanism plays a vital role in reducing background clutter and emphasizing tumor boundaries. This aligns with the findings in method. txt, where this module was introduced as a lightweight yet highly effective refinement step for anomaly localization and consistency. On the other hand, the Causal-Aware Refinement, responsible for boundary preservation and fine-level reconstruction, proves essential for segmentation precision. When excluded, both recall and F1 scores experience consistent declines-for instance, in NIH ChestX-ray14, the F1 score drops from 90.07% to 88.51%, and the AUC falls from 93.15% to 92.02%. This demonstrates that simply generating coarse masks is insufficient, and a dedicated boundary-aware structure enhances the output granularity necessary for clinical reliability.

The full model achieves top performance across all datasets and metrics, confirming the synergistic value of its integrated modules. Its architecture balances semantic abstraction with spatial detail, and consistent gains across diverse modalities-2D X-rays, 3D MRIs, and WSIs-highlight strong generalizability. These findings affirm that each component contributes meaningfully, making the framework both modular and interpretable, with clear potential for real-world clinical deployment.

# 5 Conclusions and future work

In this study, we aim to enhance the reliability and interpretability of anomaly detection in clinical settings by leveraging the power of multimodal foundation models. Traditional statistical and deep learning models, though widely used, often lack the capacity to fully capture the nuanced temporal, categorical, and semantic relationships present in medical records. To address these limitations, we propose a novel multimodal framework that combines three complementary modules. A symbolic abstraction mechanism encodes multimodal patient records into mathematically formalized representations. We introduce PathoGraph, a graph-based neural network that constructs a dynamic, symptom-centered latent space, enabling structured disentanglement of clinical variables over time. Third, the Knowledge-Guided Refinement (KGR) module integrates medical ontologies like SNOMED CT and ICD-10 via uncertainty-aware attention mechanisms and differentiable constraints. These components maintain semantic interpretability and align with medical reasoning processes. Empirical validation across real-world EHR and diagnostic datasets shows superior performance in identifying complex anomalies such as unusual combid trajectories and treatment deviations, with marked gains in robustness and transparency over baseline models.

Despite these promising results, two primary limitations remain. The framework's reliance on curated domain ontologies may limit scalability or adaptability in under-resourced clinical contexts where structured knowledge bases are incomplete or evolving. Model generalization under extreme distribution shifts-such as those caused by pandemics or rare

disease outliers-still poses a challenge, particularly when labeled data is scarce or inconsistent. Future work will explore the integration of self-supervised pretraining with broader clinical corpora and adaptive ontology expansion, aiming to enhance zero-shot adaptability and reduce domain dependency. Our study sets a foundation for explainable, multimodal AI systems in medicine, with a clear path toward broader real-world deployment.

Despite the promising results, our proposed framework has several limitations. First, the model's performance is sensitive to the quality and coverage of the external medical ontologies (e.g., SNOMED CT, ICD-10). Incomplete or outdated ontological structures may propagate semantic errors into the latent space. Second, while symbolic constraints improve interpretability, they may limit model flexibility in highly heterogeneous or emergent clinical domains such as rare diseases or novel pandemic conditions. Third, the integration of multi-modal data assumes availability of both text and imaging inputs; in cases where one modality is missing or highly noisy, the system's robustness may be reduced. In terms of computation, although the model is optimized for modular efficiency, the use of graph-based encoders and attention mechanisms does result in moderate resource requirements during training and inference. These constraints may affect scalability in low-resource clinical environments. Finally, failure cases were observed in scenarios involving ambiguous temporal sequences or overlapping symptom clusters, where disentangled representations may become less distinguishable. Future work will explore adaptive regularization, knowledge base expansion, and model compression to address these limitations and enhance deployability across broader 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.

# Ethics statement

Ethical approval was not required for the study involving humans in accordance with the local legislation and institutional requirements. Written informed consent to participate in this study was not required from the participants or the participants'; legal guardians/next of kin in accordance with the national legislation and the institutional requirements.

# Author contributions

ZeT: Conceptualization, methodology, software, validation, formal analysis, investigation, data curation, writing – original draft, Supervision, Project administration, Resources, Visualization, Writing – review and editing. ZoT: Data curation, Conceptualization, Formal analysis, Investigation, Funding acquisition, Software, Writing – original draft, Writing – review

and editing. JW; writing – review and editing, visualization, supervision, funding acquisition.

## 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.

## 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

Alhaddad, A. Y., Aly, H., Gad, H., Al-Ali, A., Sadasivuni, K. K., Cabibihan, J.-J., et al. (2022). Sense and learn: recent advances in wearable sensing and machine learning for blood glucose monitoring and trend-detection. *Front. Bioeng. Biotechnol.* 10, 876672. doi:10.3389/fbioe.2022.876672

Arora, R., Saini, I., and Sood, N. (2023). "Dual-encoder framework for multi-abnormality classification of vindr-cxr dataset," in *2023 third international conference on secure cyber computing and communication (ICSCCC)* (IEEE), 207–211.

Bayane, I., Leander, J., and Karoumi, R. (2025). Taxonomic framework for neural network-based anomaly detection in bridge monitoring. *Automation Constr.* 173, 106113. doi:10.1016/j.autcon.2025.106113

Bergmann, P., Batzner, K., Fauser, M., Sattlegger, D., and Steger, C. (2021). The mvtec anomaly detection dataset: a comprehensive real-world dataset for unsupervised anomaly detection. *Int. J. Comput. Vis.* 129, 1038–1059. doi:10.1007/s11263-020-01400-4

Cao, Y., Xu, X., Sun, C., Cheng, Y., Du, Z., Gao, L., et al. (2023). Personalizing vision-language models with hybrid prompts for zero-shot anomaly detection. *IEEE Trans. Cybern.* 55, 1917–1929. doi:10.1109/tcyb.2025.3536165

Cao, Y., Zhang, J., Frittoli, L., Cheng, Y., Shen, W., and Boracchi, G. (2024). "Adaclip: adapting clip with hybrid learnable prompts for zero-shot anomaly detection," in *European conference on computer vision* (Springer), 55–72.

Chen, W., Yi, Z., Lim, L. J. R., Lim, R. Q. R., Zhang, A., Qian, Z., et al. (2024). Deep learning and remote photoplethysmography powered advancements in contactless physiological measurement. *Front. Bioeng. Biotechnol.* 12, 1420100. doi:10.3389/fbioe.2024.1420100

Defard, T., Setkov, A., Loesch, A., and Audigier, R. (2020). Padim: a patch distribution modeling framework for anomaly detection and localization. *ICPR Work.* Available online at: https://link.springer.com/chapter/10.1007/978-3-030-68799-1_35.

DeMedeiros, K., Hendawi, A. M., and Alvarez, M. (2023). "A survey of ai-based anomaly detection in iot and sensor networks," in Italian National Conference on Sensors.

Deng, A., and Hooi, B. (2021). "Graph neural network-based anomaly detection in multivariate time series," in AAAI Conference on Artificial Intelligence, February 2-9, 2021.

Deng, H., and Li, X. (2022). *Anomaly detection via reverse distillation from one-class embedding*. Computer Vision and Pattern Recognition. Available online at: http://openaccess.thecvf.com/content/CVPR2022/html/Deng_Anomaly_Detection_via_Reverse_Distillation_From_One-Class_Embedding_CVPR_2022_paper.html.

Gudovskiy, D. A., Ishizaka, S., and Kozuka, K. (2021). "Cflow-ad: real-time unsupervised anomaly detection with localization via conditional normalizing flows," in IEEE Workshop/Winter Conference on Applications of Computer Vision, Waikoloa, HI, USA, 03-08 January 2022.

Gui, G., Gao, B.-B., Liu, J., Wang, C., and Wu, Y. (2024). "Few-shot anomaly-driven generation for anomaly classification and segmentation," in *European conference on computer vision* (Springer), 210–226.

Guo, Y., and Huang, W. (2025). "Llava-next-med: medical multimodal large language model," in *2025 asia-europe conference on cybersecurity, internet of things and soft computing (CITSC)* (IEEE), 474–477.

Hallinan, J. T. P. D., Feng, M., Ng, D., Sia, S. Y., Tiong, V. T. Y., Jagmohan, P., et al. (2022). Detection of pneumothorax with deep learning models: learning from radiologist labels vs natural language processing model generated labels. *Acad. Radiol.* 29, 1350–1358. doi:10.1016/j.acra.2021.09.013

Han, S., Hu, X., Huang, H., Jiang, M., and Zhao, Y. (2022). Adbench: anomaly detection benchmark. *Neural Inf. Process. Syst.* Available online at: https://proceedings.neurips.cc/paper_files/paper/2022/hash/cf93972b116ca5268827d575f2cc226b-Abstract-Datasets_and_Benchmarks.html.

Jiang, J., Zhu, J., Bilal, M., Cui, Y., Kumar, N., Dou, R., et al. (2023). Masked swin transformer unet for industrial anomaly detection. *IEEE Trans. Industrial Inf.* 19, 2200–2209. doi:10.1109/tii.2022.3199228

Lee, M., Park, S., Kim, H., Yoon, M., Lee, J., Choi, J. W., et al. (2024). "Spade: sparse pillar-based 3d object detection accelerator for autonomous driving," in *2024 IEEE international symposium on high-performance computer architecture (HPCA)* (IEEE), 454–467.

Li, C.-L., Sohn, K., Yoon, J., and Pfister, T. (2021). Cutpaste: self-supervised learning for anomaly detection and localization. *Comput. Vis. Pattern Recognit.* Available online at: http://openaccess.thecvf.com/content/CVPR2021/html/Li_CutPaste_Self-Supervised_Learning_for_Anomaly_Detection_and_Localization_CVPR_2021_paper.html.

Liang, H., Song, H., Zhang, S., and Bu, Y. (2025). Highway spillage detection using an improved stpm anomaly detection network from a surveillance perspective. *Appl. Intell.* 55, 7–22. doi:10.1007/s10489-024-06066-w

Liu, J., Xie, G., Wang, J., Li, S., Wang, C., Zheng, F., et al. (2023a). Deep industrial image anomaly detection: a survey. *Mach. Intell. Res.* 21, 104–135. doi:10.1007/s11633-023-1459-z

Liu, X., Tian, M., Zhu, Q., Wang, Y., Huo, H., Chen, T., et al. (2025). Selective single-bacterium analysis and motion tracking based on conductive bulk-surface imprinting. *Anal. Chem.* 97, 8915–8922. doi:10.1021/acs.analchem.5c00198

Liu, Y., Li, Z., Pan, S., Gong, C., Zhou, C., and Karypis, G. (2021). Anomaly detection on attributed networks via contrastive self-supervised learning. *IEEE Trans. Neural Netw. Learn. Syst.* 33, 2378–2392. doi:10.1109/tnnls.2021.3068344

Liu, Z., Zhou, Y., Xu, Y., and Wang, Z. (2023b). Simplenet: a simple network for image anomaly detection and localization. *Comput. Vis. Pattern Recognit.*, 20402–20411. doi:10.1109/cvpr52729.2023.01954

20252025okLet me transcribe.

Makrogiannis, S., Annasamudram, N., and Biswas, T. (2022a). "Tidaqunet: tissue identification and quantification network for mid-thigh ct segmentation," in *Medical imaging 2022: image processing* (SPIE), 12032, 109–827. doi:10.1117/12.2613304

Makrogiannis, S., Okorie, A., Di Iorio, A., Bandinelli, S., and Ferrucci, L. (2022b). Multi-atlas segmentation and quantification of muscle, bone and subcutaneous adipose tissue in the lower leg using peripheral quantitative computed tomography. *Front. Physiology* 13, 951368. doi:10.3389/fphys.2022.951368

Makrogiannis, S., Zheng, K., and Harris, C. (2021). Discriminative localized sparse approximations for mass characterization in mammograms. *Front. Oncol.* 11, 725320. doi:10.3389/fonc.2021.725320

Mishra, P., Verk, R., Fornasier, D., Piciarelli, C., and Foresti, G. (2021). Vt-adl: a vision transformer network for image anomaly detection and localization. *Int. Symposium Industrial Electron.*, 01–06. doi:10.1109/isie45552.2021.9576231

Murakami, N., Hiramatsu, N., Kobayashi, H., Akizuki, S., and Hashimoto, M. (2024). "Dc-padim: fast anomaly detection based on effective feature dimension-compression," in *2024 twelfth international symposium on computing and networking (CANDAR)* (IEEE), 190–195.

Park, H., Noh, J., and Ham, B. (2020). Learning memory-guided normality for anomaly detection. *Comput. Vis. Pattern Recognit.* Available online at: http://openaccess.thecvf.com/content_CVPR_2020/html/Park_Learning_Memory-Guided_Normality_for_Anomaly_Detection_CVPR_2020_paper.html.

Roth, K., Pemula, L., Zepeda, J., Scholkopf, B., Brox, T., and Gehler, P. (2021). *Towards total recall in industrial anomaly detection*. Computer Vision and Pattern Recognition. Available online at: http://openaccess.thecvf.com/content/CVPR2022/html/Roth_Towards_Total_Recall_in_Industrial_Anomaly_Detection_CVPR_2022_paper.html.

Su, B., Zhou, Z., and Chen, H. (2023). Pvel-ad: a large-scale open-world dataset for photovoltaic cell anomaly detection. *IEEE Trans. Industrial Inf.* 19, 404–413. doi:10.1109/tii.2022.3162846

Sun, Z., Wang, J., and Li, Y. (2024). Ramfae: a novel unsupervised visual anomaly detection method based on autoencoder. *Int. J. Mach. Learn. Cybern.* 15, 355–369. doi:10.1007/s13042-023-01913-7

Tang, J., Li, J., Gao, Z.-C., and Li, J. (2022). "Rethinking graph neural networks for anomaly detection," in International Conference on Machine Learning, Baltimore, Maryland USA, July 17-23, 2022.

Tian, Y., Pang, G., Chen, Y., Singh, R., Verjans, J., and Carneiro, G. (2021). "Weakly-supervised video anomaly detection with robust temporal feature magnitude learning," in IEEE International Conference on Computer Vision, Montreal, QC, Canada, 10-17 October 2021.

Tien, T. D., Nguyen, A. T., Tran, N. H., Huy, T. D., Duong, S. T., Tr, C. D., et al. (2023). *Revisiting reverse distillation for anomaly detection*. Computer Vision and Pattern Recognition. Available online at: https://openaccess.thecvf.com/content/CVPR2023/html/Tien_Revisiting_Reverse_Distillation_for_Anomaly_Detection_CVPR_2023_paper.html?ref=https://githubhelp.com.

Tuli, S., Casale, G., and Jennings, N. (2022). "Tranad: deep transformer networks for anomaly detection in multivariate time series data," in *Proceedings of the VLDB endowment*.

Wang, D., Zhuang, L., Gao, L., Sun, X., Huang, M., and Plaza, A. (2023a). Bocknet: blind-block reconstruction network with a guard window for hyperspectral anomaly detection. *IEEE Trans. Geoscience Remote Sens.* 61, 1–16. doi:10.1109/tgrs.2023.3335484

Wang, Y., Peng, J., Zhang, J., Yi, R., Wang, Y., and Wang, C. (2023b). *Multimodal industrial anomaly detection via hybrid fusion*. Computer Vision and Pattern Recognition. Available online at: https://openaccess.thecvf.com/content/CVPR2023/html/Wang_Multimodal_Industrial_Anomaly_Detection_via_Hybrid_Fusion_CVPR_2023_paper.html?ref=https://githubhelp.com.

Wijerathna, V., Raveen, H., Abeygunawardhana, S., and Ambegoda, T. D. (2022). "Chest x-ray caption generation with chexnet," in *2022 moratuwa engineering research conference (MERCon)* (IEEE), 1–6.

Wyatt, J., Leach, A., Schmon, S. M., and Willcocks, C. G. (2022). "Anoddpm: anomaly detection with denoising diffusion probabilistic models using simplex noise," in 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), New Orleans, LA, USA, 19-20 June 2022, 649–655. doi:10.1109/cvprw56347.2022.00080

Xu, J., Wu, H., Wang, J., and Long, M. (2021). "Anomaly transformer: time series anomaly detection with association discrepancy," in International Conference on Learning Representations, Vienna, Austria, 4 May 2021.

Yang, Y., Zhang, C., Zhou, T., Wen, Q., and Sun, L. (2023). "Dcdetector: dual attention contrastive representation learning for time series anomaly detection," in *Knowledge discovery and data mining*.

You, Z., Cui, L., Shen, Y., Yang, K., Lu, X., Zheng, Y., et al. (2022). A unified model for multi-class anomaly detection. *Neural Inf. Process. Syst.* Available online at: https://proceedings.neurips.cc/paper_files/paper/2022/hash/1d774c112926348c3e25ea47d87c835b-Abstract-Conference.html.

Zavrtanik, V., Kristan, M., and Skočaj, D. (2021). "DrÃ†m - a discriminatively trained reconstruction embedding for surface anomaly detection," in IEEE International Conference on Computer Vision, Montreal, QC, Canada, 10-17 October 2021, 8310–8319. doi:10.1109/iccv48922.2021.00822

Zhang, Y., Dong, Q., Zhao, X., Sun, Y., Lin, X., Zhang, X., et al. (2024). Honeycomb-like biomimetic scaffold by functionalized antibacterial hydrogel and biodegradable porous mg alloy for osteochondral regeneration. *Front. Bioeng. Biotechnol.* 12, 1417742. doi:10.3389/fbioe.2024.1417742

Zhou, Q., Pang, G., Tian, Y., He, S., and Chen, J. (2023). "Anomalyclip: object-agnostic prompt learning for zero-shot anomaly detection," in International Conference on Learning Representations, Kigali, Rwanda, May 1-5, 2023.

Zou, Y., Jeong, J., Pemula, L., Zhang, D., and Dabeer, O. (2022). "Spot-the-difference self-supervised pre-training for anomaly detection and segmentation," in European Conference on Computer Vision, Tel Aviv, Israel, October 23–27, 2022, 392–408. doi:10.1007/978-3-031-20056-4_23