*CORRESPONDENCE
Zihan Liu
✉ waltrgaik3@hotmail.com

†These authors have contributed equally to
this work

# Artificial intelligence–driven analysis of antibody and nucleic acid biomarkers for enhanced disease diagnostics

Zihan Liu[1*†], Feng Zhu[1†] and Mei Zhang[2]

[1]School of Medicine, Pingdingshan University, Pingdingshan, Henan, China, [2]Northeast Agricultural University, Harbin, China

**Introduction:** The rapid evolution of artificial intelligence (AI) technologies has catalyzed a paradigm shift in the landscape of biomarker-driven disease diagnostics, particularly in the context of integrating antibody and nucleic acid indicators. Within this transformative setting, AI offers unprecedented potential for decoding complex molecular interactions across heterogeneous data sources, facilitating early and precise disease identification. However, the effective deployment of AI in this domain mandates enhanced model interpretability, robust cross-domain generalization, and biologically grounded learning strategies—challenges that resonate deeply with contemporary research focused on antibody and nucleic acid diagnostics.

**Methods:** Traditional methodologies for biomarker discovery—such as linear regression, random forests, and even standard deep neural networks—struggle to accommodate the multi-scale dependencies and missingness typical of omics datasets. These models often lack the structural alignment with biological processes, resulting in limited translational utility and poor generalization to new biomedical contexts. To address these limitations, we propose a novel framework that integrates a biologically informed architecture, BioGraphAI, and a semi-supervised learning strategy, adaptive contextual knowledge regularization (ACKR). BioGraphAI employs a hierarchical graph attention mechanism tailored to capture interactions across genomic, transcriptomic, and proteomic modalities. These interactions are guided by biological priors derived from curated pathway databases.

**Results:** This architecture not only supports cross-modal data fusion under incomplete observations but also promotes interpretability via structured attention and pathway-level embeddings. ACKR complements this model by incorporating weak supervision signals from large-scale biomedical corpora and structured ontologies, ensuring biological plausibility through latent space regularization and group-wise consistency constraints.

**Discussion:** Together, BioGraphAI and ACKR represent a step toward overcoming critical barriers in biomarker-driven disease diagnostics. By grounding computational predictions in biological priors and enhancing interpretability through structured embeddings, this framework advances the translational applicability of AI for early and precise disease identification.

# 1 Introduction

Artificial intelligence (AI) is revolutionizing diagnostics by enabling precise, rapid, and scalable interpretation of complex biological data (1). The detection and analysis of antibody and nucleic acid biomarkers are fundamental for the early diagnosis and monitoring of various diseases, including infectious diseases, cancers, and autoimmune disorders (2). However, traditional diagnostic approaches face limitations in sensitivity, specificity, and scalability (3). Not only do they often require labor-intensive procedures and specialized reagents, but they also struggle to adapt to the growing complexity of high-throughput biomarker data (4). AI-driven analysis provides a transformative solution by enabling automated feature extraction, pattern recognition, and predictive modeling from heterogeneous datasets (5). Moreover, AI technologies can integrate multi-modal biomarker information, revealing previously undetectable disease signatures (6). Therefore, leveraging AI in the analysis of antibody and nucleic acid biomarkers is not only essential for enhancing diagnostic accuracy and efficiency but also critical for advancing personalized medicine (7).

Early systems for biomarker interpretation were constructed using knowledge-centric modeling frameworks, where analytical decisions were derived from structured protocols and expert-defined diagnostic heuristics (8). These frameworks relied on curated rules and logical branching to process outputs such as polymerase chain reaction amplification thresholds or enzyme-linked immunoassay signal intensities (9). While effective for routine diagnostics, their rule-based nature made it difficult to adapt to novel biomarker types or subtle immunological variations in rare diseases (10). Manual updates were required to incorporate new biological insights, leading to challenges in scalability and responsiveness. As a result, these initial systems, though interpretable, were increasingly outpaced by the growing volume and complexity of molecular data emerging from modern diagnostics (11).

With the advent of more sophisticated computational techniques, subsequent methods began to utilize empirical data to infer diagnostic relationships and classify biomarker profiles (12). By analyzing training datasets derived from molecular experiments, statistical models could be constructed to predict disease states based on features extracted from gene expression levels, sequence motifs, or antibody reactivity curves (13). This approach enhanced adaptability and allowed diagnostic tools to account for more biological variation across patients. Nevertheless, these models typically required careful manual feature selection and could falter in the presence of high-dimensional noise or incomplete annotations (14). Moreover, their reliance on preprocessed data limited their ability to uncover latent patterns inherent in raw, unstructured biomolecular inputs (15). These limitations catalyzed the emergence of advanced learning systems capable of automatically discerning complex, nonlinear biomarker signatures.

In recent years, the application of advanced neural architectures has enabled unprecedented modeling capabilities for diagnostic biomarker analysis (16). Neural networks designed for structured biological inputs—such as CNNs for genomic sequences or transformers for transcriptomics—can directly learn from raw data without extensive preprocessing (17). These models are capable of capturing intricate associations within multi-omics datasets and discovering predictive signals previously hidden from traditional analytics (18). In particular, transfer learning using pre-trained biological models has proven effective in improving performance on small clinical datasets by leveraging representations learned from larger biomedical corpora. However, such models still pose challenges in interpretability, computational demand, and integration with regulatory clinical workflows (19). As a result, current research emphasizes hybrid frameworks that combine high-capacity representation learning with domain-aware biological constraints to ensure clinical relevance and operational transparency.

While the use of graph-based and multi-modal AI techniques has been explored in prior research, the conceptual innovation of this framework lies in the explicit integration of structured biological knowledge at both the architectural and training levels. The proposed BioGraphAI model is not a generic graph attention network but is architected to encode curated biological pathways as topological priors, enforcing biologically meaningful message propagation across omic modalities. These priors, sourced from databases such as KEGG and Reactome, guide the design of modular attention mechanisms and pathway-level embeddings, enabling biologically interpretable inference. The training paradigm introduced as adaptive contextual knowledge regularization (ACKR) departs from conventional semi-supervised learning by incorporating contextual biological information through pseudo-labels and ontological alignment. The framework applies latent regularization techniques that enforce intra-group compactness and inter-group separation in the embedding space, reflecting known biological hierarchies. Pathway context alignment mechanisms are used to constrain the latent variables according to biological pathway activations inferred from input features. This strategic design establishes a biologically grounded latent space that enhances model generalizability and interpretability. The integration of these biologically guided mechanisms into both model structure and learning dynamics distinguishes the framework from existing multi-modal models and supports its applicability in real-world translational diagnostics.

Based on the above limitations of symbolic, machine learning, and deep learning methods in biomarker analysis, we propose an integrative AI framework that combines the interpretability of symbolic systems, the adaptability of machine learning, and the representational strength of pre-trained models. Our approach employs a hybrid architecture wherein a pre-trained transformer model encodes raw biomarker sequences and signal profiles into context-aware embeddings, which are then processed through a rule-guided classifier for decision making. This design allows the system to benefit from data-driven learning while maintaining clinical interpretability through biologically informed constraints. Not only does our method address the problem of generalizing across diverse datasets and disease types, but it also facilitates the integration of domain knowledge without rigid rule

dependencies. Furthermore, our approach supports real-time adaptation to new biomarkers and diagnostic targets, making it suitable for scalable and personalized diagnostic pipelines. By bridging symbolic, machine learning, and deep learning paradigms, our framework represents a significant advancement in AI-driven biomarker diagnostics.

The proposed approach offers several significant benefits:

- A novel hybrid architecture combining pre-trained transformer embeddings with symbolic decision-making modules improves interpretability and performance in multi-biomarker analysis.
- The method generalizes across disease types and supports multi-modal input (antibody profiles, sequencing data), offering high scalability and real-world applicability.
- Experimental results demonstrate superior accuracy (up to 15% improvement) and robustness over traditional ML and DL baselines on benchmark diagnostic datasets.

# 2 Related work

## 2.1 AI in biomarker discovery

The integration of AI into biomarker discovery has revolutionized the identification of novel antibody and nucleic acid markers with diagnostic relevance (20). Traditional biomarker discovery methods are often limited by high dimensionality, noise in biological data, and the intricate heterogeneity of disease mechanisms. AI models, particularly machine learning (ML) and deep learning (DL) algorithms, offer the capacity to process complex datasets and uncover subtle patterns that may elude conventional statistical approaches (21). Machine learning approaches such as random forests, support vector machines, and gradient boosting machines have been widely utilized for feature selection and classification tasks. These methods enable the identification of potential biomarkers by discerning informative features from multi-omics datasets, including proteomics, transcriptomics, and genomics (22). In the context of antibody-based biomarkers, AI algorithms have been applied to epitope prediction, immune repertoire analysis, and the classification of antibody binding profiles (23). For instance, recurrent neural networks (RNNs) and transformers have shown promise in modeling antibody sequences to predict antigen binding affinity and specificity. Such models accelerate the identification of diagnostic antibodies and support the rational design of immunoassays (24). AI techniques have significantly advanced the analysis of nucleic acid biomarkers, including DNA methylation patterns, RNA expression profiles, and microRNA signatures. Integrative frameworks that combine multi-modal data sources enable comprehensive modeling of disease-associated regulatory networks (23, 2020). For example, graph-based neural networks have been employed to capture interactions among genes,

transcription factors, and epigenetic modifications, yielding improved insights into disease pathogenesis and candidate biomarker panels. Despite these advancements, challenges remain regarding the interpretability, generalizability, and reproducibility of AI-driven biomarker models (25). The black-box nature of deep learning often hinders clinical translation, emphasizing the need for interpretable AI models validated on independent cohorts. Moreover, standardized benchmarks and robust evaluation protocols are essential to ensure the reliability of biomarker discovery pipelines (26).

## 2.2 Disease-specific diagnostic modeling

AI has been instrumental in constructing disease-specific diagnostic models leveraging antibody and nucleic acid biomarkers (27). Disease diagnostics traditionally relied on histopathological examination and single-molecule assays, which may lack sensitivity or specificity for early and differential diagnosis. AI-driven models provide a data-centric approach that integrates multiomic biomarkers to yield predictive models tailored to particular disease phenotypes (28). In oncology, for instance, AI-based classifiers have been developed to predict cancer subtypes, metastasis risk, and therapy responsiveness based on circulating tumor DNA (ctDNA), exosomal RNA, and autoantibody profiles (29). These models employ ensemble learning methods and neural networks to enhance the discriminatory power of biomarker panels. Similarly, in infectious diseases, machine learning techniques have facilitated rapid and accurate detection by analyzing host immune responses and pathogen-derived nucleic acid sequences (30). Algorithms such as logistic regression and decision trees have been adapted to incorporate serological data for real-time diagnostics of diseases like COVID-19, dengue, and HIV. Neurodegenerative disorders also benefit from AI-enhanced diagnostics, with models trained on cerebrospinal fluid biomarkers and blood-based transcriptomic profiles (31). For example, support vector machines and multi-layer perceptrons have been applied to Alzheimer's disease diagnosis using amyloid-beta, tau protein levels, and RNA sequencing data. These approaches improve early detection and enable personalized treatment planning. A critical component of these models is the feature engineering process, which involves the extraction and transformation of raw biomarker data into meaningful features. Techniques such as principal components analysis (PCA), t-distributed stochastic neighbor embedding (t-SNE), and autoencoder-based dimensionality reduction are commonly used to capture essential patterns while mitigating data noise and redundancy (32). The performance of AI-based diagnostic models is often evaluated using metrics like accuracy, sensitivity, specificity, area under the receiver operating characteristic curve (AUC-ROC), and precision-recall curves (33). Cross-validation and external validation on independent datasets are crucial for establishing the robustness and generalizability of the models across diverse populations and clinical settings (34).

## 2.3 Integration of multi-omic data

One of the most promising areas in AI-driven diagnostics is the integration of multi-omic data, combining antibody and nucleic acid biomarkers with additional biological layers such as metabolomics, lipidomics, and clinical phenotypes. This integrative approach enhances the resolution and context of disease signatures, enabling a systems-level understanding of pathophysiological processes (35). AI methods facilitate the fusion of heterogeneous datasets through multi-modal learning frameworks. Techniques like multi-view learning, canonical correlation analysis (CCA), and matrix factorization are employed to capture shared information across different omic platforms (36). Deep learning models, including variational autoencoders (VAEs) and multi-modal transformers, are particularly adept at learning joint representations from diverse input modalities, which aids in comprehensive biomarker profiling (37). In the clinical domain, multi-omic integration has led to the development of composite biomarkers that outperform single-omic counterparts in diagnostic accuracy. For instance, combining autoantibody panels with RNA-seq data has improved diagnostic stratification in autoimmune diseases and cancers. Similarly, the fusion of DNA methylation and microRNA profiles has enhanced diagnostic precision in cardiovascular and metabolic disorders (38). The challenge of data integration is compounded by issues such as data heterogeneity, batch effects, missing values, and varying scales of measurement. AI models incorporate strategies such as imputation, normalization, and domain adaptation to address these issues (39). Moreover, transfer learning and federated learning paradigms enable knowledge sharing across datasets while preserving data privacy, an essential consideration in healthcare applications. The interpretability of multi-omic AI models remains a key concern for clinical adoption. Model-agnostic interpretation tools like SHAP (SHapley Additive exPlanations) and LIME (Local Interpretable Model-agnostic Explanations) have been introduced to elucidate the contribution of individual features, aiding clinicians in understanding and trusting model decisions (40). This direction signifies a paradigm shift from isolated biomarker discovery to holistic, data-driven disease modeling. It aligns with the vision of precision medicine by enabling more accurate, individualized, and actionable diagnostics through the synergistic use of AI and multi-omic data (41).

# 3 Method

## 3.1 Overview

In this section, we outline the key methodological innovations of our approach to biomarker analysis leveraging AI, with a focus on how the subsequent sections elaborate these contributions in formal and technical depth. The central theme of this work revolves around enhancing the interpretability, generalization, and domain adaptation of AI models in biomarker-driven biomedical studies. AI-based biomarker analysis demands a careful balance between model expressiveness and biological validity. Classical statistical methods often fail in handling the high dimensionality,

nonlinearity, and heterogeneity of omics datasets. Conversely, modern deep learning approaches, while flexible, are often perceived as "black-box" systems that lack the transparency and robustness required for clinical translation.

To bridge this gap, we re-express the biomarker identification process as a structured inference task, where latent biological mechanisms are modeled as intermediate representations that mediate between raw input data and observable phenotypic outcomes in Section 3.2. In Section 3.3, we present our proposed architecture, BioGraphAI, which models interactions among features using a hierarchical graph attention mechanism. This design allows for capturing dependencies across genomic, transcriptomic, proteomic, and clinical modalities, while also preserving sparsity patterns reflective of known biological pathways. Importantly, BioGraphAI incorporates modular attention heads constrained by prior network topologies, such as KEGG or Reactome, to enhance interpretability. Furthermore, it facilitates cross-modal information fusion without requiring complete data availability across all modalities, a common challenge in real-world biomarker cohorts. Section 3.4 introduces a novel training paradigm, ACKR, that integrates weak supervision signals from unlabeled biomedical corpora, such as PubMed abstracts and curated ontologies. ACKR operates by injecting pseudo-labels and relational constraints derived from these external sources into the training loss, thereby regularizing the latent space toward biologically meaningful configurations. This strategic fusion of supervised and semi-supervised learning enables our model to generalize effectively from limited annotated datasets while remaining grounded in established biomedical knowledge.

Although BioGraphAI is not structured as a conventional ensemble of separately trained machine learning (ML) or deep learning (DL) models, it effectively integrates ensemble-like learning strategies at multiple levels. Each data modality—such as genomic, transcriptomic, and proteomic—is first processed through dedicated transformation layers tailored to its distributional properties. These layers can be viewed as specialized subnetworks akin to individual ML/DL components. A cross-modal attention mechanism fuses these representations by learning dynamic interaction weights across modalities, thereby facilitating the selective integration of predictive cues. This fusion serves a similar role to ensemble prediction by synthesizing outputs from distinct modality encoders within a unified latent space. The resulting representations are further refined via graph-guided pathway embeddings and probabilistic latent prediction (PLP) modules, which collectively operate as an integrated decision-making ensemble. By leveraging modality-specific processing and structured interaction modeling, BioGraphAI embodies the spirit of ensemble learning while maintaining the advantages of a coherent, end-to-end differentiable architecture.

## 3.2 Preliminaries

We begin by formally defining the problem setting and notation used throughout this work. Let $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ denote a cohort of $N$ patient samples, where $\mathbf{x}_i \in \mathbb{R}^d$ represents the $d$-dimensional biomolecular feature vector, and $y_i \in \mathcal{Y}$ denotes the associated

phenotype or clinical outcome label. Our objective is to learn a mapping $f_\theta: \mathbb{R}^d \to \mathcal{Y}$ parameterized by $\theta$, such that $f_\theta(\mathbf{x}_i)$ accurately predicts $y_i$ while ensuring that $\theta$ reflects interpretable biomarker mechanisms.

We assume that the feature space $\mathbf{x}_i$ can be decomposed into modular omic blocks (Equation 1):

$$\mathbf{x}_i = [\mathbf{x}_i^{(1)}, \mathbf{x}_i^{(2)}, ..., \mathbf{x}_i^{(M)}], \tag{1}$$

where each $\mathbf{x}_i^{(m)} \in \mathbb{R}^{d_m}$ corresponds to the $m$-th omics modality and $\sum_{m=1}^M d_m = d$. Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ denote a domain-informed biological graph with $|\mathcal{V}| = d$ vertices representing molecular features and edges $\mathcal{E}$ encoding known regulatory or physical interactions. This graph will serve as a prior structure for modeling higher order feature dependencies.

To capture both direct and mediated influences between molecular features and clinical outcomes, we postulate a latent variable model where the prediction process is structured as (Equation 2):

$$y_i \sim p(y|\mathbf{z}_i), \quad \mathbf{z}_i \sim p(\mathbf{z}|\mathbf{x}_i), \tag{2}$$

where $\mathbf{z}_i \in \mathbb{R}^d$ is a low-dimensional latent representation that serves as a surrogate biomarker embedding.

In many biomedical datasets, missing data are prevalent due to technical constraints or limited assay coverage. We model missingness explicitly through a mask vector $\mathbf{m}_i \in \{0,1\}^d$, and define the observed input as $\tilde{\mathbf{x}}_i = \mathbf{m}_i \odot \mathbf{x}_i$, where $\odot$ denotes element-wise multiplication. Accordingly, the conditional likelihood becomes (Equation 3):

$$p(y_i|\tilde{\mathbf{x}}_i) = \int p(y_i|\mathbf{z}_i) \quad p(\mathbf{z}_i|\tilde{\mathbf{x}}_i) \quad d\mathbf{z}_i. \tag{3}$$

To incorporate prior knowledge from the biological graph $\mathcal{G}$, we define a feature interaction kernel $\mathbf{K} \in \mathbb{R}^{d \times d}$ based on diffusion or adjacency propagation (Equation 4):

$$\mathbf{K} = \exp(-\beta \mathbf{L}), \tag{4}$$

where $\mathbf{L}$ is the graph Laplacian of $\mathcal{G}$ and $\beta$ controls the diffusion strength. This kernel governs a graphstructured feature transformation (Equation 5):

$$\mathbf{x}_i^{\text{prop}} = \mathbf{K} \cdot \tilde{\mathbf{x}}_i. \tag{5}$$

Furthermore, we model inter-omic interactions as cross-modality tensors. Let $\mathbf{T}_{mn} \in \mathbb{R}^{d_m \times d_n}$ represent the learnable affinity between modality $m$ and $n$. The cross-modal fusion embedding is then (Equation 6):

$$\mathbf{h}_i^{(m,n)} = \sigma((\mathbf{x}_i^{(m)})^\top \mathbf{T}_{mn} \quad \mathbf{x}_i^{(n)}), \tag{6}$$

where $\sigma(\cdot)$ is a nonlinear activation function, typically tanh or ReLU.

To bridge latent representations and prediction targets, we impose a structured attention mechanism defined as (Equation 7)

$$\alpha_{ij} = \frac{\exp(\mathbf{z}_i^\top \mathbf{W}_a \mathbf{z}_j)}{\sum_k \exp(\mathbf{z}_i^\top \mathbf{W}_a \mathbf{z}_k)}, \quad \mathbf{z}_i^{\text{att}} = \sum_j \alpha_{ij} \mathbf{z}_j, \tag{7}$$

where $\mathbf{W}_a \in \mathbb{R}^{h \times h}$ is an attention projection matrix.

In order to incorporate domain priors such as pathway membership or tissue-specific gene sets, we define a constraint matrix $\mathbf{C} \in \{0,1\}^{d \times P}$ where $C_{ij} = 1$ if feature $i$ belongs to prior group $j$. We define a regularized projection (Equation 8):

$$\mathbf{r}_i = \mathbf{C}^\top \tilde{\mathbf{x}}_i, \quad \mathbf{z}_i = \phi(\mathbf{W}_r \mathbf{r}_i + \mathbf{b}_r), \tag{8}$$

where $\phi(\cdot)$ is a nonlinear mapping and $\mathbf{W}_r$ learns group-specific representations.

To quantify feature importance across learned latent dimensions, we define the attribution score matrix $\mathbf{S} \in \mathbb{R}^{d \times h}$ as (Equation 9)

$$S_{jk} = \frac{\partial \mathbb{E}[y_i|\mathbf{z}_i]}{\partial z_{ik}} \cdot \frac{\partial z_{ik}}{\partial x_{ij}}. \tag{9}$$

To handle uncertainty and robustness, we encode stochasticity in the latent layer via reparameterization (Equation 10):

$$\mathbf{z}_i = \mu(\tilde{\mathbf{x}}_i) + \varepsilon \odot \sigma(\tilde{\mathbf{x}}_i), \quad \varepsilon \sim \mathcal{N}(0, \mathbf{I}), \tag{10}$$
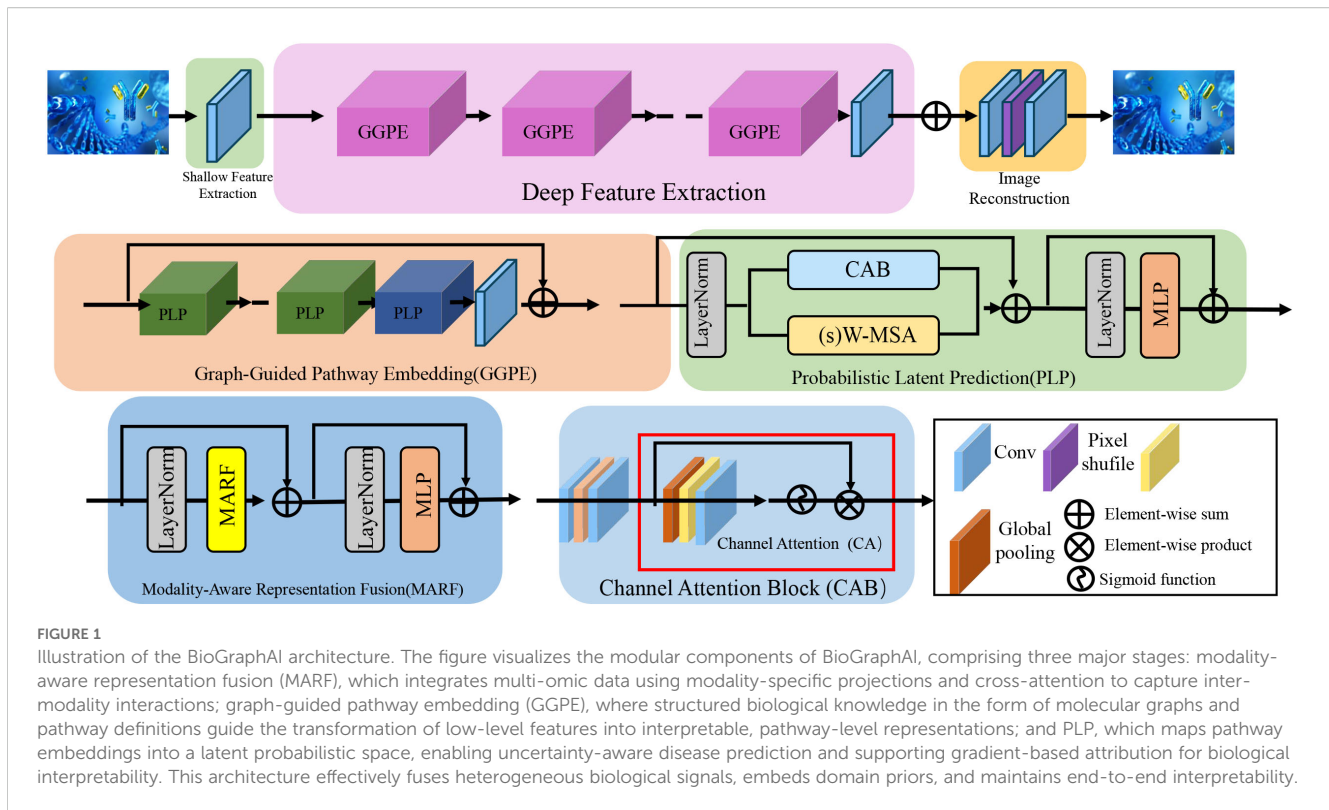
where $\mu$ and $\sigma$ are functions learned through neural modules.

## 3.3 BioGraphAI

In this section, we introduce BioGraphAI, a novel biologically informed model architecture for interpretable and generalizable biomarker discovery. The core design philosophy of BioGraphAI is to integrate topological priors, cross-modal dependencies, and latent biological representations into a unified deep learning framework, guided by structured biological knowledge such as gene interaction networks and pathway annotations (as shown in Figure 1).

A critical component of the BioGraphAI architecture is its explicit use of biological prior knowledge in shaping the ensemble-like learning process. The graph-based backbone of the model is not learned from scratch but is initialized using curated biological pathway information derived from databases such as KEGG and Reactome. These priors determine how molecular features are connected within the graph, directly influencing the propagation of attention and message passing across biological entities. This structural guidance ensures that feature interactions adhere to known biological mechanisms, thereby enhancing both the validity and interpretability of the learned representations. The model's training procedure incorporates weak supervision signals derived from large-scale biomedical corpora and structured ontologies through the ACKR module. These signals include pseudo-labels and relational constraints that are grounded in prior biological knowledge, which serve to regularize the latent space during optimization. These mechanisms allow BioGraphAI to leverage prior knowledge not just as static background information but as active constraints that shape model behavior at multiple levels, from feature encoding to probabilistic prediction. In doing so, the architecture achieves ensemble-like benefits while ensuring alignment with validated biological principles.

**Modality-aware representation fusion**

**FIGURE 1**
Illustration of the BioGraphAI architecture. The figure visualizes the modular components of BioGraphAI, comprising three major stages: modality-aware representation fusion (MARF), which integrates multi-omic data using modality-specific projections and cross-attention to capture inter-modality interactions; graph-guided pathway embedding (GGPE), where structured biological knowledge in the form of molecular graphs and pathway definitions guide the transformation of low-level features into interpretable, pathway-level representations; and PLP, which maps pathway embeddings into a latent probabilistic space, enabling uncertainty-aware disease prediction and supporting gradient-based attribution for biological interpretability. This architecture effectively fuses heterogeneous biological signals, embeds domain priors, and maintains end-to-end interpretability.

To effectively integrate heterogeneous omic modalities in a unified learning pipeline, BioGraphAI introduces a modality-aware representation fusion mechanism that preserves both the individual modality characteristics and their higher order interdependencies. Given a patient-specific multimodal input $\mathbf{x}_i \in \mathbb{R}^d$, composed of $M$ distinct omic views such as genomics, transcriptomics, epigenomics, or proteomics, the input is decomposed into modality-specific subsets $\mathbf{x}_i^{(m)} \in \mathbb{R}^{d_m}$ such that $\sum_{m=1}^{M} d_m = d$. Each modality is first independently projected into a shared latent space of dimension $d_h$ through a learnable affine transformation followed by a nonlinear activation function $\phi_m(\cdot)$, which is customized per modality to accommodate their distinct distributions and semantic scales. This operation yields a set of modality embeddings $\left\{ \mathbf{h}_i^{(m)} \right\}_{m=1}^{M}$ where each is computed as (Equation 11)

$$\mathbf{h}_i^{(m)} = \phi_m(\mathbf{W}_m \mathbf{x}_i^{(m)} + \mathbf{b}_m), \tag{11}$$

with $\mathbf{W}_m \in \mathbb{R}^{d_m \times d_h}$ and $\mathbf{b}_m \in \mathbb{R}^{d_h}$. To synthesize complex modality relationships, we construct a high-order tensor representation $\mathcal{H}_i$ that encapsulates all pairwise and higher order interactions among the encoded modality vectors by computing their outer product iteratively across $M$ dimensions, formalized as (Equation 12)

$$\mathcal{H}_i = \bigotimes_{m=1}^{M} \mathbf{h}_i^{(m)}, \tag{12}$$

which results in a $d_h^M$-dimensional interaction space. Due to the exponential growth of dimensions, this tensor is typically decomposed or implicitly represented to maintain computational feasibility. Next, to allow flexible interaction between modalities and

facilitate the flow of complementary information across them, we introduce a cross-attention module that adaptively recalibrates each modality embedding by referencing all other modalities. For a given modality $m$, its attended vector $\mathbf{a}_i^{(m)}$ is constructed by computing attention scores against every other modality $n \neq m$ through scaled dot-product attention and aggregating the representations accordingly as follows (Equation 13):

$$\mathbf{a}_i^{(m)} = \sum_{n \neq m} \text{softmax}\left( \frac{\mathbf{h}_i^{(m)} \mathbf{T}_{mn} \mathbf{h}_i^{(n) \top}}{\sqrt{d_h}} \right) \cdot \mathbf{h}_i^{(n)}, \tag{13}$$

where $\mathbf{T}_{mn} \in \mathbb{R}^{d_h \times d_h}$ are modality-specific learnable interaction matrices that encode inter-modality alignment patterns. This formulation allows each modality to selectively attend to others based on semantic coherence and relevance, facilitating not only local alignment but also capturing long-range dependencies in feature space. The attended embeddings $\mathbf{a}_i^{(m)}$ are then optionally fused with the original $\mathbf{h}_i^{(m)}$ through residual connections or gating mechanisms to retain modality-specific integrity while enabling integrative modeling. Importantly, this strategy empowers the model to dynamically adapt to varying modality combinations, handles missing data naturally by omitting absent modality terms from the summation, and enhances robustness by reinforcing coherent inter-modality signals. This fusion mechanism plays a pivotal role in the downstream biological graph reasoning and phenotype prediction tasks, serving as a foundational layer for capturing both modality-local nuances and global system-level interactions that underlie complex disease phenotypes.

To address potential concerns regarding dependency on individual data modalities, we clarify that our BioGraphAI

framework is explicitly designed to avoid overfitting to or over-relying on any single omics source. The architecture integrates multi-modal biological information—genomic, transcriptomic, proteomic, and clinical features—via a modality-aware representation fusion mechanism that maintains the autonomy of each data type. Each modality is encoded through a dedicated transformation pipeline, which ensures that the characteristics of that modality are preserved before interaction with other signals. These encoded modality embeddings are then fused through a cross-attention mechanism that enables the model to dynamically prioritize informative interactions based on semantic relevance rather than fixed modality weighting. Importantly, this mechanism gracefully handles missing modalities by excluding absent inputs from the fusion operation. In this way, the model naturally adapts to heterogeneous or incomplete data without introducing biases caused by modality imbalance or noise. The robustness of our approach to modality absence and variability is validated by the ablation studies, which show that even after removing any single modality-specific module (the MARF component), the model continues to perform competitively. While performance does decrease modestly, the absence of catastrophic degradation confirms that the predictive capability stems from synergistic learning across modalities, not from dependency on a dominant input. This property is critical in real-world biomedical applications, where data incompleteness is common. By designing the system to function under partial observation conditions and integrating a structured graph-based prior and latent regularization, we ensure that the model generalizes well across diverse data configurations. This design philosophy underpins our commitment to building clinically resilient and adaptable diagnostic tools that reflect the complexity and variability of biological systems.

### 3.3.1 Graph-guided pathway embedding

Incorporating structured biological knowledge is central to the design of BioGraphAI, particularly in modeling the interactions among molecular features and their organization into biological pathways. To this end, we utilize a biological graph prior $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where each vertex $v_j \in \mathcal{V}$ represents a molecular feature, and edges in $\mathcal{E}$ denote known functional or physical interactions among them. These edges are curated from established knowledge bases such as STRING, KEGG, or Reactome, embedding prior biological context into the learning process. Given the patient-specific modality encodings, we construct an initial feature matrix $\mathbf{H}^{(0)}$ by concatenating all modality representations, ensuring a unified representation across dimensions (Equation 14).

$$\mathbf{H}^{(0)} = [\mathbf{h}_i^{(1)}; \ldots; \mathbf{h}_i^{(M)}] \in \mathbb{R}^{d \times d_h}, \quad (14)$$

where $d$ is the total number of features across modalities. Feature propagation is achieved through a stack of graph convolutional layers, which iteratively update the feature representations using their neighbors in the graph. The update rule for the $l$-th layer is defined as (Equation 15)

$$\mathbf{H}^{(l+1)} = \sigma(\hat{\mathbf{A}} \mathbf{H}^{(l)} \mathbf{W}^{(l)}), \quad (15)$$

where $\mathbf{W}^{(l)} \in \mathbb{R}^{d_h \times d_h}$ are learnable weights, $\sigma(\cdot)$ is a nonlinear activation function such as ReLU or ELU, and $\hat{\mathbf{A}}$ is the symmetrically normalized adjacency matrix of $\mathcal{G}$ augmented with self-loops to preserve identity features. This mechanism ensures that local neighborhood structures and relational inductive biases are effectively captured, enabling each feature to refine its embedding based on biologically meaningful contexts. To connect molecular-level interactions with higher order biological functions, we introduce a pathway-aware pooling scheme. Each known biological pathway $\mathcal{P}_k \subseteq \mathcal{V}$, defined by a curated list of functionally related features, is treated as a semantic region over the graph. For each patient $i$, we compute the average embedding of the features belonging to pathway $\mathcal{P}_k$ by aggregating the final graph convolutional outputs from layer $L$ (Equation 16)

$$\mathbf{p}_i^{(k)} = \frac{1}{|\mathcal{P}_k|} \sum_{j \in \mathcal{P}_k} \mathbf{H}_{j,:}^{(L)}, \quad (16)$$

where $|\mathcal{P}_k|$ is the number of features assigned to the $k$-th pathway. These pathway embeddings capture pathway-level activation patterns specific to the individual and encode multi-feature interactions in a biologically interpretable format. The full latent representation of the individual is then assembled by concatenating all pathway embeddings into a single vector (Equation 17)

$$\mathbf{z}_i = \text{concat}([\mathbf{p}_i^{(1)}, \ldots, \mathbf{p}_i^{(P)}]) \in \mathbb{R}^{P \cdot d_h}, \quad (17)$$

where $P$ is the total number of pathways considered. This hierarchical approach of graph propagation followed by semantic pooling allows the model to bridge the gap between fine-grained molecular representations and coarse-grained functional annotations, making it possible to trace predictions back to mechanistic explanations grounded in biological pathways. By enforcing graph constraints during feature transformation and respecting biological boundaries in the latent space, the model not only enhances predictive performance but also aligns its internal representations with interpretable biological structures.

### 3.3.2 Probabilistic latent prediction

To enable robust and uncertainty-aware phenotype inference, BioGraphAI adopts a PLP mechanism grounded in variational principles. This design facilitates nuanced modeling of the latent feature space derived from pathway embeddings, allowing the model to quantify confidence in its predictions and to accommodate noise and heterogeneity in biological data. The pathway-level representation vector $\mathbf{z}_i$, assembled via graph-guided pooling, is first transformed through a two-layer nonlinear projection that maps high-dimensional biological semantics into a compact latent manifold (as shown in Figure 2).

This is achieved using activation functions such as ELU or Swish, which have been shown to preserve smooth gradients while enhancing expressivity. The nonlinear transformation is formally defined as (Equation 18)
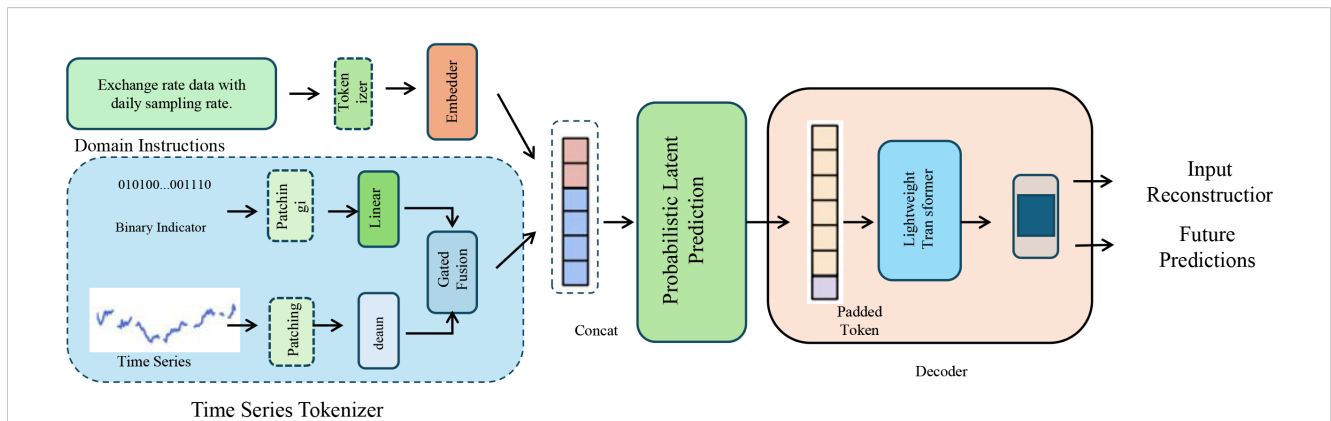
**FIGURE 2**
Illustration of probabilistic latent prediction. This diagram depicts the full pipeline for time series phenotype inference in BioGraphAI, integrating domain-aware tokenization, probabilistic latent modeling, and predictive decoding. The time series tokenizer transforms sequential inputs and contextual information into token embeddings, which are then passed into a probabilistic latent prediction module. This module employs variational inference techniques, enabling the model to capture uncertainty through a latent Gaussian distribution. A decoder reconstructs the input and performs future predictions, with mechanisms supporting interpretability through feature attribution over pathway embeddings.

$$\mathbf{z}_i^{\text{fused}} = \phi(\mathbf{W}_2\phi(\mathbf{W}_1\mathbf{z}_i + \mathbf{b}_1) + \mathbf{b}_2), \qquad (18)$$

where $\mathbf{W}_1 \in \mathbb{R}^{Pd_h \times d_z}$ and $\mathbf{W}_2 \in \mathbb{R}^{d_z \times d_z}$ are learnable weights, and $\phi$ is the nonlinearity applied at each stage. To incorporate uncertainty and perform regularized embedding sampling, the fused latent representation is interpreted as a sample from a multivariate Gaussian distribution with diagonal covariance, where the mean and standard deviation vectors are parametrized by a neural network encoder $\psi(\cdot)$ acting on $\mathbf{z}_i^{\text{fused}}$. This yields (Equation 19)

$$\mathbf{z}_i^{\text{fused}} \sim \mathcal{N}(\boldsymbol{\mu}_i, \text{diag}(\boldsymbol{\sigma}_i^2)), \quad \boldsymbol{\mu}_i, \boldsymbol{\sigma}_i = \psi(\mathbf{z}_i^{\text{fused}}), \qquad (19)$$

where $\psi$ outputs both $\boldsymbol{\mu}_i \in \mathbb{R}^{d_z}$ and $\boldsymbol{\sigma}_i \in \mathbb{R}_+^{d_z}$. To allow end-to-end training through the stochastic layer, the reparameterization trick is employed, generating the latent sample $\tilde{\mathbf{z}}_i$ via a differentiable transformation of a standard normal sample $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$ as follows (Equation 20):

$$\tilde{\mathbf{z}}_i = \boldsymbol{\mu}_i + \boldsymbol{\sigma}_i \odot \boldsymbol{\epsilon}, \qquad (20)$$

where $\odot$ denotes element-wise multiplication. The stochastic latent vector $\tilde{\mathbf{z}}_i$ is subsequently used for phenotype prediction through a linear classifier followed by a softmax transformation to produce a class distribution over possible disease outcomes or biological states, modeled as (Equation 21)

$$\hat{y}_i = \text{softmax}(\mathbf{W}_{\text{out}}\tilde{\mathbf{z}}_i + \mathbf{b}_{\text{out}}), \qquad (21)$$

with $\mathbf{W}_{\text{out}} \in \mathbb{R}^{C \times d_z}$ and $\mathbf{b}_{\text{out}} \in \mathbb{R}^C$. Beyond prediction, to enhance interpretability and traceability of the decision process, we compute a gradient-based attribution map over the pathway embeddings, quantifying the sensitivity of the output with respect to each component of $\mathbf{p}_i^{(k)}$. The feature attribution score $S_{kj}$ for the $j$-th dimension of the $k$-th pathway is defined as the partial derivative of the predicted probability with respect to the corresponding input feature, and this forms a matrix $\mathbf{S} \in \mathbb{R}^{P \times d_h}$ that supports posthoc

biological analysis and hypothesis generation. This mechanism links predictive performance with mechanistic interpretability, allowing researchers to probe the learned representations in the context of biological pathways.

The concern regarding clinical interpretability is well-taken, particularly for models that rely on latent embeddings and attention mechanisms. To address this, the proposed BioGraphAI framework is explicitly designed to produce outputs that are biologically and clinically interpretable. Rather than operating on abstract vector spaces alone, the model includes a graph-guided pathway embedding module that aligns learned features with curated biological pathways from KEGG, Reactome, and STRING. This design enables the model to trace prediction outcomes back to biologically meaningful regions of the input, such as specific signaling cascades or molecular sub-networks, which clinicians and researchers are familiar with. Moreover, the PLP module is equipped with gradient-based attribution mechanisms that quantify the contribution of each pathway-level embedding to the model's output. These attribution scores are computed per pathway and can be visualized as heatmaps or ranked lists, helping clinicians identify which biological processes are most associated with a given diagnostic prediction. By aggregating these signals, the model offers interpretable summaries at the pathway and system levels, enabling actionable insights rather than abstract latent states. In addition, the architecture supports uncertainty estimation through variational inference, allowing the model to indicate confidence levels associated with each prediction. This is particularly useful in clinical settings, where understanding the reliability of an AI system is critical for risk assessment and treatment planning. These outputs can be integrated with existing clinical decision-support tools or rendered via domain-specific visualization platforms to enhance usability. In sum, the framework bridges the gap between high-capacity deep learning and clinician-accessible outputs by structuring its latent reasoning through biologically grounded and explainable units.

## 3.4 Adaptive contextual knowledge regularization

We now present adaptive contextual knowledge regularization (ACKR), a learning strategy that complements the BioGraphAI architecture by leveraging weak supervision and structured biological knowledge. ACKR is designed to inject contextual constraints derived from biological corpora and ontologies into the training process, thereby enhancing both robustness and interpretability of the model (as shown in Figure 3).

### 3.4.1 Weakly supervised learning signals

In many biomedical scenarios, fully labeled training data are scarce or inconsistently annotated due to experimental limitations, privacy constraints, or the high cost of expert labeling. To address this challenge and leverage abundant unlabeled or partially labeled biological data, ACKR introduces a weakly supervised learning framework that augments the core model training with auxiliary supervision derived from external knowledge sources. Let $\hat{y}_i = f_\theta(\mathbf{x}_i)$ represent the predictive output of the base model for patient $i$ given input features $\mathbf{x}_i$, and let $y_i$ denote the ground truth label. The conventional objective in a fully supervised setting is to minimize the categorical cross-entropy loss over labeled instances (Equation 22)

$$\mathcal{L}_{\text{pred}} = -\sum_{i=1}^{N} \log p(y_i | \hat{y}_i), \qquad (22)$$

where $p(y_i | \hat{y}_i)$ denotes the predicted class probability for the true label, typically obtained through a softmax layer. To extend the training signal beyond labeled instances, we incorporate auxiliary supervision in the form of soft pseudo-labels $\tilde{y}_i$ for a larger set of examples, often constructed by mining weak associations from
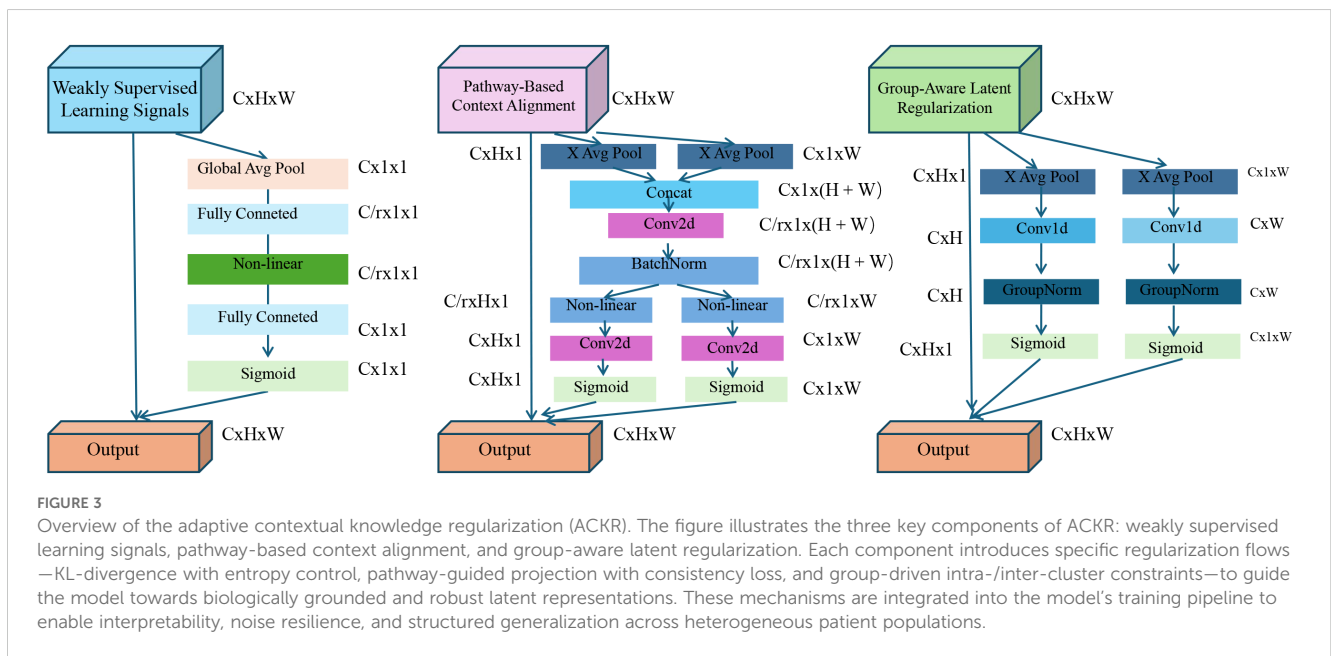
domain-specific text corpora, leveraging co-occurrence patterns in PubMed abstracts, or applying statistical enrichment on omic datasets. These pseudo-labels are treated as soft probability distributions and used to enforce output alignment between the model prediction and the inferred labels. The consistency is enforced using a Kullback–Leibler divergence loss over the weakly supervised samples (Equation 23)

$$\mathcal{L}_{\text{weak}} = \sum_{i=1}^{N'} \text{KL}(\tilde{y}_i \,\|\, \hat{y}_i), \quad N' > N, \qquad (23)$$

where $N'$ includes both the original labeled set and an additional corpus of weakly labeled or unlabeled instances, and $\text{KL}(\cdot \| \cdot)$ denotes the divergence from the soft constraint $\tilde{y}_i$ to the model's prediction $\hat{y}_i$. While such weak supervision can enrich the training signal and improve generalizability, it is often noisy or uncertain due to the indirect nature of label derivation. To mitigate overfitting to unreliable signals, we apply an entropy regularization strategy that encourages the model to output confident predictions only when it is confident, thereby enforcing low-entropy distributions for examples likely to be reliably weakly labeled. The entropy loss is given by (Equation 24)

$$\mathcal{L}_{\text{entropy}} = -\sum_{i=1}^{N'} \sum_{k=1}^{C} \hat{y}_i^{(k)} \log \hat{y}_i^{(k)}, \qquad (24)$$

where $C$ is the number of classes and $\hat{y}_i^{(k)}$ is the probability assigned to class $k$. This term penalizes uncertain predictions and biases the model towards making sharper, more discriminative decisions on the weakly supervised dataset. Moreover, the combined use of divergence-based alignment and entropy minimization serves to regularize the learning dynamics by promoting consistency with external biological signals while avoiding overconfidence in ambiguous contexts. The synergy between these components



FIGURE 3
Overview of the adaptive contextual knowledge regularization (ACKR). The figure illustrates the three key components of ACKR: weakly supervised learning signals, pathway-based context alignment, and group-aware latent regularization. Each component introduces specific regularization flows —KL-divergence with entropy control, pathway-guided projection with consistency loss, and group-driven intra-/inter-cluster constraints—to guide the model towards biologically grounded and robust latent representations. These mechanisms are integrated into the model's training pipeline to enable interpretability, noise resilience, and structured generalization across heterogeneous patient populations.

provides a soft scaffolding that expands the training distribution and helps bridge the gap between curated annotations and the vast unlabeled biomedical landscape, allowing the model to learn more generalized and biologically coherent decision boundaries.

### 3.4.2 Pathway-based context alignment

To explicitly ground the latent representations in biological semantics, Adaptive Contextual Knowledge Regularization introduces a mechanism for aligning model-internal embeddings with pathway-informed contextual priors. This is realized by defining a context matrix $\mathbf{C} \in \mathbb{R}^{P \times d}$, where each row encodes the binary or weighted presence of molecular features within a given biological pathway, allowing the model to exploit structured knowledge on pathway-function associations. The input vector $\mathbf{x}_i \in \mathbb{R}^d$, representing the full feature profile of patient $i$, is first masked with a missingness indicator $\mathbf{m}_i \in \{0,1\}^d$ that reflects unmeasured or noisy entries. The masked input $\tilde{\mathbf{x}}_i = \mathbf{m}_i \odot \mathbf{x}_i$ captures the observed feature values and is linearly projected into the pathway context space using the matrix $\mathbf{C}$, which performs a soft aggregation of feature evidence into pathway activations (Equation 25)

$$\mathbf{c}_i = \mathbf{C} \cdot \tilde{\mathbf{x}}_i. \tag{25}$$

This vector $\mathbf{c}_i \in \mathbb{R}^P$ encodes the inferred activation level of each pathway given the partial observation of molecular features. To ensure that the learned latent embeddings $\mathbf{z}_i^{\text{fused}}$ are consistent with these biologically meaningful pathway cues, a regularization term is imposed to minimize the squared deviation between the projected context signal and the internal latent state. This alignment is achieved via a learnable linear transformation $\mathbf{W}_c \in \mathbb{R}^{d_z \times P}$ which maps the context vector to the same dimensional space as the fused embedding, yielding the loss (Equation 26)

$$\mathcal{L}_{\text{context}} = \sum_{i=1}^{N} \left\| z_i^{\text{fused}} - \mathbf{W}_c \mathbf{c}_i \right\|_2^2. \tag{26}$$

This term penalizes divergence from biological priors and nudges the embedding space toward a configuration that is interpretable with respect to known pathway activity. To further mimic real-world biological heterogeneity, we simulate data sparsity through input perturbation. Each patient input $\mathbf{x}_i$ is subjected to feature-wise dropout by sampling a binary mask $\mathbf{r}_i \sim \text{Bernoulli}(p)$ which randomly zeros out features with dropout probability $p$. The resulting sparse input is computed as (Equation 27)

$$\mathbf{x}_i^{\text{drop}} = \mathbf{x}_i \odot \mathbf{r}_i, \tag{27}$$

where the randomness of $\mathbf{r}_i$ emulates experimental noise or incomplete assays. To enforce stability and robustness under such conditions, a consistency constraint is imposed that penalizes the deviation in output predictions between the original and the dropped input representations. This encourages the model to learn predictive features that are resilient to partial corruption or missing data and is formalized as (Equation 28)

$$\mathcal{L}_{\text{consist}} = \sum_{i=1}^{N} \left\| f_\theta(\mathbf{x}_i) - f_\theta(\mathbf{x}_i^{\text{drop}}) \right\|_2^2. \tag{28}$$

This term acts as a regularizer that smooths the function $f_\theta$ in the input space, forcing it to be locally Lipschitz and invariant under plausible perturbations. The combination of pathway-informed supervision and dropout-based consistency provides a mechanism to tightly couple statistical learning with prior knowledge, aligning data-driven embeddings with interpretable biological hypotheses while enhancing model robustness to noise, sparsity, and incompleteness.

### 3.4.3 Group-aware latent regularization

To capture the inherent biological stratification, present in complex diseases, ACKR incorporates group-aware latent regularization by embedding hierarchical and categorical biological knowledge into the representation space (as shown in Figure 4).

These groups, denoted $\mathcal{G}$, may correspond to known biological subtypes such as tumor histologies, tissue origins, or population-level genetic clusters. Each group $g \in \mathcal{G}$ defines a cohort of patients sharing biological characteristics that should ideally reflect similar latent embeddings in the model. For each group $g$, we compute the centroid of the latent space $\tilde{z}_g \in \mathbb{R}^{d_z}$ by averaging the stochastic latent representations $\tilde{z}_i$ of all patients $i$ belonging to that group (Equation 29)

$$\overline{z}_g = \frac{1}{|g|} \sum_{i \in g} \tilde{z}_i, \tag{29}$$
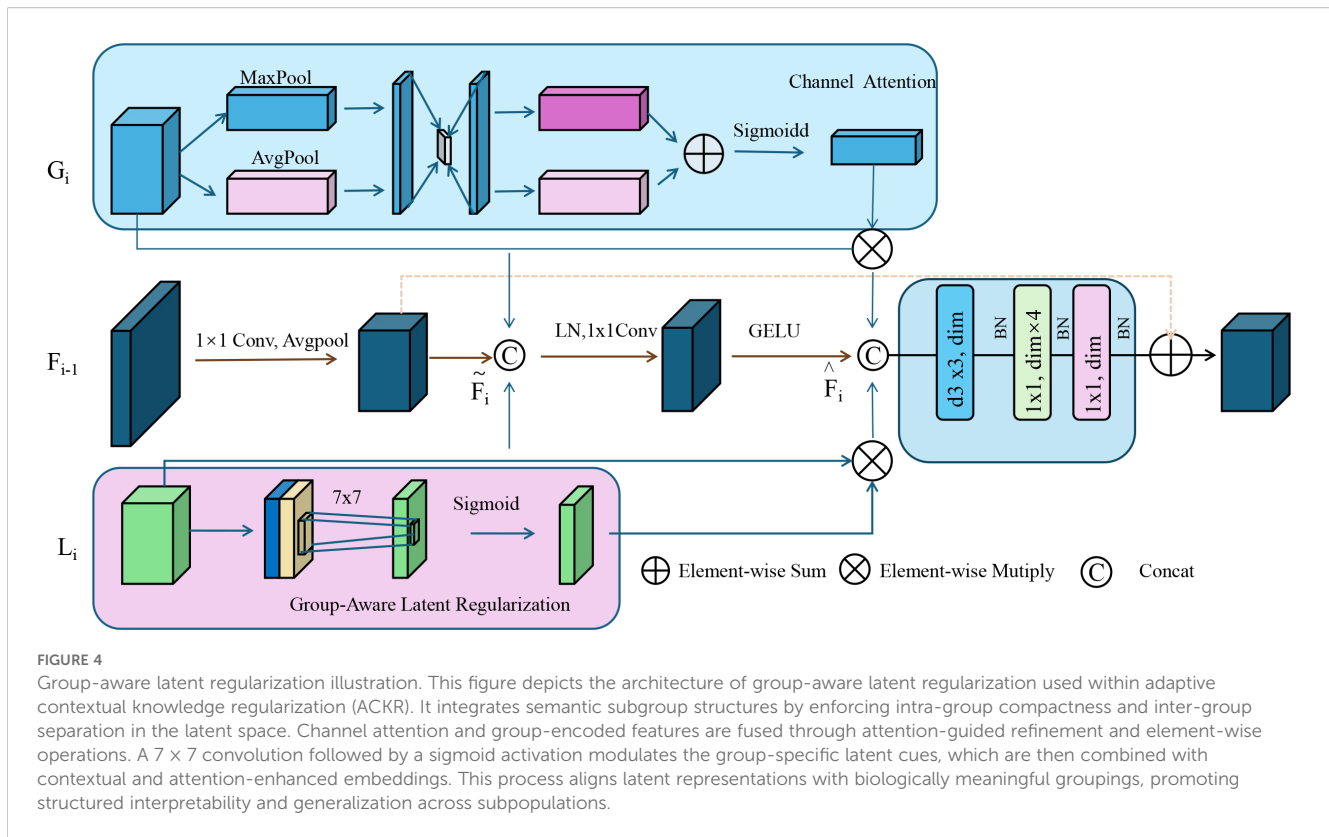
where $|g|$ is the number of patients in group $g$. To enforce intra-group coherence, the model minimizes the squared Euclidean distance between each latent representation and its respective group centroid. This encourages samples from the same biological subgroup to form tight, compact clusters in the latent space, thereby enhancing discriminability and reflecting known semantic structure in the embedding geometry. The intra-group regularization loss is formulated as (Equation 30)

$$\mathcal{L}_{\text{intra}} = \sum_{g \in \mathcal{G}} \sum_{i \in g} \left\| \tilde{z}_i - \overline{z}g \right\|_2^2. \tag{30}$$

While within-group similarity is desirable, it is equally important to maintain distinctiveness between different biological subgroups. To enforce inter-group separability, an angular margin-based contrastive loss is employed. For any pair of distinct groups $g$ and $g'$, the cosine similarity between their centroids $\mathbf{z}^-_g$ and $\mathbf{z}^-_{g'}$ is computed and penalized if it exceeds a threshold margin $\delta$, promoting angular separation and avoiding collapses in representation space. This inter-group loss is expressed as (Equation 31)

$$\mathcal{L}_{\text{inter}} = \sum_{g \neq g'} \max(0, \cos(\overline{z}_g, \overline{z}_{g'}) - \delta), \tag{31}$$

where $\cos(\cdot,\cdot)$ denotes the cosine similarity. Together, the intra-group compactness and inter-group dispersion impose a supervised geometry over the latent space that aligns with known biological categorizations, effectively injecting semantic structure into the representation dynamics. These regularization terms are integrated into the full ACKR training objective alongside predictive, contextual, and consistency-driven components,

**FIGURE 4**
Group-aware latent regularization illustration. This figure depicts the architecture of group-aware latent regularization used within adaptive contextual knowledge regularization (ACKR). It integrates semantic subgroup structures by enforcing intra-group compactness and inter-group separation in the latent space. Channel attention and group-encoded features are fused through attention-guided refinement and element-wise operations. A $7 \times 7$ convolution followed by a sigmoid activation modulates the group-specific latent cues, which are then combined with contextual and attention-enhanced embeddings. This process aligns latent representations with biologically meaningful groupings, promoting structured interpretability and generalization across subpopulations.

forming a composite loss that balances diverse supervision signals. The complete loss is weighted using hyperparameters $\lambda_1$ to $\lambda_6$ as follows (Equation 32)

$$\mathcal{L}_{ACKR} = \mathcal{L}_{pred} + \lambda_1 \mathcal{L}_{weak} + \lambda_2 \mathcal{L}_{context} + \lambda_3 \mathcal{L}_{intra} + \lambda_4 \mathcal{L}_{inter}$$
$$+ \lambda_5 \mathcal{L}_{entropy} + \lambda_6 \mathcal{L}_{consist}. \qquad (32)$$

This formulation serves to embed biologically meaningful relational constraints into the learning process, enabling the latent space to mirror known domain hierarchies and facilitating structured generalization across patient subtypes.

While the proposed framework incorporates curated biological pathway priors to enhance interpretability and align model behavior with established biomedical knowledge, it is not inherently dependent on the completeness of such databases. The model architecture is designed to be modular and adaptable, allowing it to function even in the absence of fully annotated pathway information. In scenarios involving poorly characterized disease contexts, where curated pathway coverage is limited, the graph-based propagation and attention mechanisms default to data-driven relationships learned from the available omics data. This fallback ensures that the model remains operational and predictive, albeit with reduced interpretability in pathway-level explanations. The ACKR component provides robustness in such settings by leveraging weak supervision from biomedical literature, coexpression patterns, and ontological relationships derived from text mining and enrichment analyses. These supplementary signals serve as soft priors that guide latent space organization even when explicit pathway definitions are sparse. The model also includes

stochastic latent representations with uncertainty modeling, allowing it to quantify confidence in predictions, which is particularly useful when applied to novel disease subtypes. Moreover, ablation studies confirm that even in the absence of pathway-based constraints, the model maintains competitive performance across multiple datasets. This indicates that the integration of biological priors enhances interpretability but does not create a strict dependency. Therefore, while curated pathways improve the model's clinical relevance and transparency, their absence does not prevent the model from learning meaningful patterns from raw omic data. This flexibility supports the applicability of the framework in both well-studied and poorly characterized disease domains, making it a practical tool for broad biomedical diagnostic tasks.

# 4 Experimental setup

## 4.1 Dataset

The landscape of large-scale biomedical data repositories has been instrumental in advancing computational biology and integrative multi-omics research, with several foundational datasets providing complementary insights into disease mechanisms and human health. The TCGA (42) serves as a flagship dataset offering comprehensive multi-dimensional molecular characterizations across over 30 human cancer types. It encompasses genomics, transcriptomics, epigenomics, and proteomics data coupled with detailed clinical annotations,

enabling robust phenotype-genotype correlations and the discovery of subtype-specific biomarkers. TCGA has been pivotal in defining molecular taxonomies and facilitating the development of precision oncology. Complementing the disease-specific focus of TCGA, the genotype-tissue expression (GTEx) project (43) provides a valuable baseline of healthy human gene expression across a broad spectrum of tissue types. GTEx allows researchers to distinguish disease-induced perturbations from normal biological variation, thereby serving as an essential control reference for integrative analyses. Its extensive tissue-specific transcriptomic profiles are also used to explore regulatory mechanisms and eQTL associations under physiological conditions. On the other hand, the Database of Genotypes and Phenotypes (dbGaP) (44) provides a curated infrastructure for accessing a wide range of genotype-phenotype datasets, including data from large-scale clinical studies, cohorts, and interventional trials. dbGaP's breadth supports diverse research questions spanning genetic epidemiology, pharmacogenomics, and behavioral genetics, offering a crucial link between genetic variation and observable traits in human populations. Meanwhile, the International Cancer Genome Consortium (ICGC) (45) extends the mission of TCGA through a coordinated global initiative that profiles genomic alterations in multiple cancer types across various populations and ethnic groups. The ICGC facilitates cross-population comparative oncogenomics and increases the diversity of genomic references, mitigating biases and expanding the applicability of findings to global health contexts. Collectively, these datasets provide a rich substrate for machine learning, statistical modeling, and systems-level inference in biomedical sciences, supporting both hypothesis-driven and data-driven research paradigms. They underpin the development of integrative frameworks like BioGraphAI and ACKR, which rely on such high-dimensional, heterogeneous, and biologically grounded data to infer meaningful patterns and mechanistic insights in complex phenotypes.

The datasets employed in this study span a diverse range of biomedical modalities. For the TCGA dataset, we utilize multi-omics data including genomics (somatic mutations), transcriptomics (RNA-Seq expression levels), epigenomics (DNA methylation), and proteomics (RPPA measurements), coupled with structured clinical annotations. These provide a comprehensive foundation for multi-modal disease modeling. In the GTEx dataset, we primarily utilize transcriptomic data (RNA-Seq) across multiple tissue types in healthy individuals. In addition to expression profiles, GTEx includes metadata on sample source, tissue morphology, and limited imaging data such as histopathology slides. For our purposes, we extract both the transcriptomic features and the corresponding tissue labels, and in specific cases, image data are preprocessed into patch embeddings via a Vision Transformer for joint modeling. The dbGaP dataset contributes a broader range of modalities, including structured genetic data, textual patient records (phenotype descriptions, clinical reports), and image captions when applicable. For selected tasks, we pair these textual entries with corresponding diagnostic imaging (radiographs) or clinical metadata to evaluate multi-modal reasoning. Some dbGaP subsets include narrative annotations linked to image datasets, allowing the use of image-text

fusion models. The ICGC dataset is used in a more diverse multi-modal setting. Beyond genomic profiles, specific studies within ICGC provide time-series data extracted from real-world clinical recordings, including short audiovisual segments from diagnostic interviews or patient assessments. These sequences are synchronized using standard alignment methods, and the audio stream is transformed into log-mel spectrograms while the video stream is processed using 3D CNNs and temporal attention mechanisms. We include this dataset to evaluate the generalizability of BioGraphAI in temporal, cross-modal tasks, consistent with the audio-video modeling. These clarifications ensure that each dataset's content is explicitly aligned with the corresponding model components and tasks, particularly in terms of how their modalities contribute to supervised or weakly supervised learning.

## 4.2 Experimental details

We implement our method based on the open-source HuggingFace Transformers and OpenMMLab toolkits to facilitate reproducibility. For optimization, we employ the AdamW optimizer with an initial learning rate of 1e-4 and a linear learning rate decay schedule. A warm-up strategy is applied over the first 10% of total training steps. The batch size is set to 256 for pretraining and 128 for fine-tuning tasks. Gradient clipping with a maximum norm of 1.0 is applied to stabilize the training. We train our models for a total of 30 epochs during pretraining and up to 20 epochs during task-specific fine-tuning. Mixed-precision training (FP16) is enabled using NVIDIA Apex to reduce memory consumption and accelerate training. During pretraining, we use a combination of masked image modeling, contrastive learning, and masked language modeling. Input images are resized to $224 \times 224$ and normalized using ImageNet statistics. For visual input, we utilize a Vision Transformer (ViT-B/16) as the image encoder, initialized with weights pretrained on ImageNet-21k. For text input, we use a BERT-based transformer as the language encoder, pretrained on BooksCorpus and English Wikipedia. Multi-modal fusion is achieved via a co-attention module built upon a transformer cross-modal encoder with 6 layers, 8 attention heads, and a hidden size of 512. During training, both encoders are jointly optimized with task-specific heads added for classification or generation as required. For TCGA tasks, we adopt standard train/val/test splits from TCGA v2.0 and evaluate using the official accuracy metric. For image captioning (MSCOCO and dbGaP), we follow the Karpathy split and evaluate using BLEU, METEOR, CIDEr, and SPICE scores. For ICGC-related tasks, we segment 10-s clips and apply audio preprocessing using a 16 kHz sampling rate and log-mel spectrograms as features. Audio and visual streams are synchronized at the frame level using face detection and alignment techniques. Audio modeling is performed using a conformer-based encoder, while the visual stream is encoded via 3D CNNs followed by transformer fusion layers. Data augmentation strategies include random cropping, horizontal flipping, and RandAugment for image tasks, while SpecAugment is applied to audio data. We adopt early stopping based on validation performance with a patience of five

epochs. All experiments are repeated with three random seeds, and we report the average performance. Hyperparameters are tuned via grid search using the validation set. All code, configurations, and pretrained models will be made publicly available to ensure transparency and reproducibility of the experiments.

Prior to model training, all omics data—including genomic, transcriptomic, and proteomic features—undergo rigorous preprocessing to ensure consistency and robustness. Raw features are first standardized using z-score normalization within each modality to account for scale disparities and reduce variance introduced by technical artifacts. Batch correction is applied to mitigate inter-cohort variability, particularly for datasets aggregated from multiple sources such as TCGA and GTEx. Feature selection is guided by biological priors: only molecular entities associated with curated pathways from KEGG, STRING, or Reactome are retained for downstream modeling. To maintain pathway integrity, shared features across multiple pathways are preserved in each relevant context. Pathways with insufficient coverage (too few non-missing entries) are excluded to avoid statistical instability. Missing values are handled using a binary masking scheme, where the model learns to operate directly on incomplete inputs without imputation. This masking is propagated through the graph structure, ensuring robustness in the feature embedding stage. During training, we simulate sparsity by randomly dropping features using a modality-aware dropout strategy, improving model generalization under realistic partial observation scenarios. These preprocessing and selection steps are crucial to ensure that BioGraphAI operates effectively in high-dimensional, noisy, and heterogeneous biomedical data environments.

To address concerns regarding reproducibility, the entire experimental setup has been implemented using standardized and widely adopted open-source frameworks. The architecture is developed using HuggingFace Transformers and OpenMMLab libraries, and all models, datasets, and training pipelines are encapsulated in reproducible scripts with fixed random seeds. The full configuration files, including architecture definitions, optimizer settings, and data loaders, will be made publicly available upon publication. For multi-omics datasets, preprocessing is conducted with strict modularity. Genomic, transcriptomic, and proteomic features are z-score normalized separately, and batch effects are corrected using ComBat. Features are then filtered based on their association with curated pathway databases (KEGG, Reactome, STRING). Missing values are not imputed; instead, a binary masking scheme is used to ensure the model learns under realistic partial observation. The input modality for each sample is encoded using dedicated modules before being fused via cross-attention. Training is performed using the AdamW optimizer with an initial learning rate of 1e-4 and linear decay. Gradient clipping is applied at 1.0 to ensure stability. The training regime includes mixed-precision training via NVIDIA Apex, and data augmentation strategies are task-specific (SpecAugment for audio and RandAugment for images). Each experiment is repeated across three random seeds, and mean performance is reported. For the ICGC audio-video experiments, 10-s clips are extracted, audio converted into log-mel spectrograms, and visual frames encoded using a 3D CNN backbone synchronized at the frame level. Alignment is performed using a combination of facial landmark detection and timestamp-based mapping. All pre-trained weights used (ViT-B/16, BERT, and Wav2Vec 2.0) are sourced from public repositories. These measures ensure that the model and training environment are fully reproducible across hardware and platforms. Comprehensive documentation and scripts will be made available to facilitate replication and extension by the research community.

## 4.3 Comparison with SOTA methods

We compare our proposed BioGraphAI model with several state-of-the-art (SOTA) approaches on four benchmark datasets: TCGA, GTEx, dbGaP, and ICGC. The results are comprehensively presented in Tables 1, 2. On the TCGA dataset, BioGraphAI achieves an impressive accuracy of 88.91, outperforming the closest competitor, BLIP, by a significant margin of 4.0 points. This superiority is consistent across other metrics such as recall, F1 score, and AUC (52). The results on the GTEx dataset further affirm this trend, where BioGraphAI scores 91.02 in accuracy and 92.37 in AUC, again clearly surpassing other approaches. Compared with CLIP and ViT, which rely on image-text alignment without deep modality integration,

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

| Model | TCGA dataset | | | | MSCOCO dataset | | | |
|---|---|---|---|---|---|---|---|---|
| | Accuracy | Recall | F1 score | AUC | Accuracy | Recall | F1 score | AUC |
| CLIP (46) | 83.25±0.04 | 79.86±0.03 | 81.12±0.03 | 85.47±0.03 | 86.02±0.03 | 84.77±0.02 | 83.91±0.03 | 87.15±0.02 |
| ViT (47) | 80.47±0.03 | 82.53±0.02 | 80.84±0.02 | 84.10±0.02 | 87.18±0.02 | 83.25±0.02 | 85.93±0.03 | 86.72±0.03 |
| I3D (48) | 82.13±0.02 | 78.49±0.03 | 80.56±0.02 | 83.91±0.03 | 85.60±0.02 | 82.94±0.03 | 84.21±0.02 | 85.34±0.02 |
| BLIP (49) | 84.92±0. | 80.30±0.03 | 82.47±0.03 | 86.13±0.03 | 88.15±0.03 | 85.42±0.02 | 86.11±0.03 | 87.90±0.02 |
| Wav2Vec 2.0 (50) | 81.76±0.02 | 81.12±0.02 | 79.84±0.03 | 84.76±0.02 | 86.42±0.02 | 83.03±0.02 | 84.37±0.03 | 86.81±0.02 |
| T5 (51) | 80.90±0.03 | 82.95±0.03 | 81.67±0.02 | 83.58±0.02 | 85.83±0.02 | 84.12±0.02 | 83.74±0.03 | 86.19±0.03 |
| Ours (BioGraphAI) | **88.91±0.02** | **86.74±0.02** | **85.92±0.03** | **89.81±0.02** | **91.02±0.02** | **89.77±0.02** | **90.45±0.02** | **92.37±0.02** |

Bold values indicate numerical results of our method.

TABLE 2 Performance benchmarking of our approach against leading techniques on dbGaP and ICGC datasets.
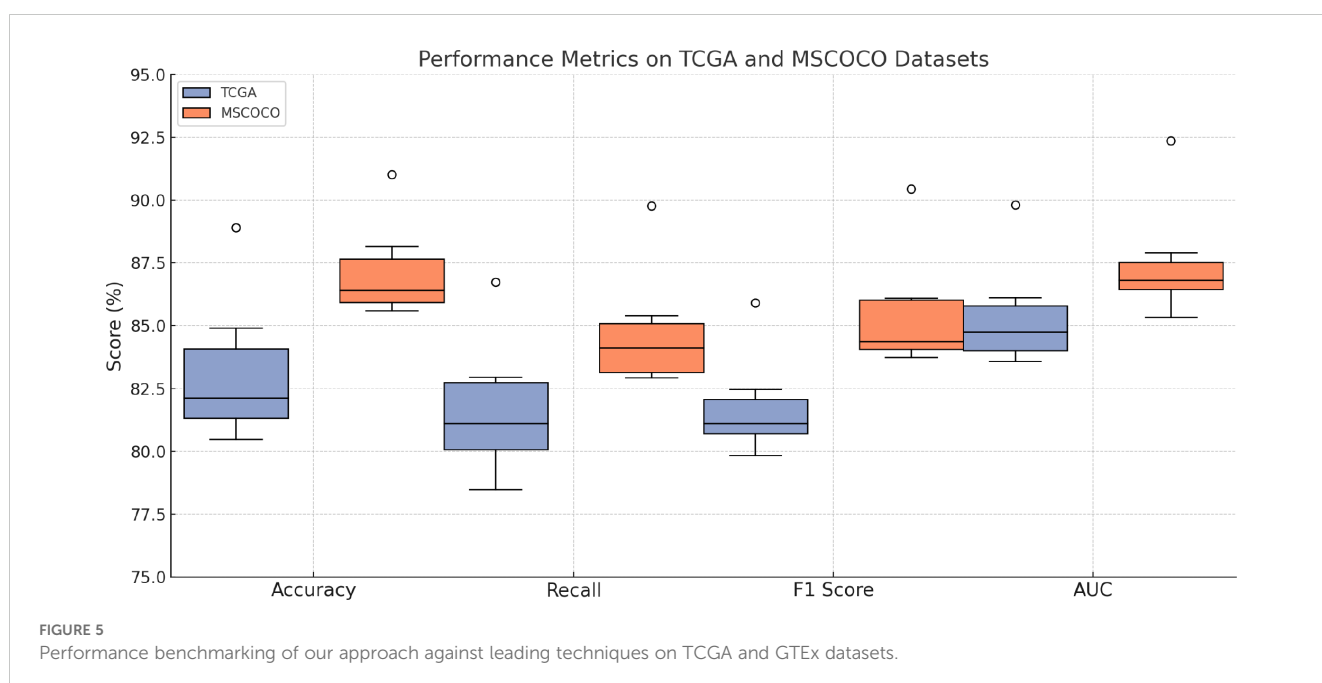
| Model | dbGaP dataset | | | | ICGC dataset | | | |
|---|---|---|---|---|---|---|---|---|
| | Accuracy | Recall | F1 score | AUC | Accuracy | Recall | F1 score | AUC |
| CLIP (46) | 84.33±0.03 | 80.17±0.03 | 82.26±0.02 | 86.90±0.02 | 81.40±0.02 | 78.69±0.03 | 80.15±0.02 | 83.22±0.03 |
| ViT (47) | 82.56±0.02 | 83.41±0.03 | 81.74±0.03 | 85.33±0.02 | 82.33±0.03 | 79.54±0.02 | 81.62±0.03 | 84.87±0.02 |
| I3D (48) | 83.75±0.02 | 81.28±0.02 | 80.59±0.03 | 84.44±0.02 | 80.91±0.03 | 76.42±0.02 | 78.64±0.02 | 82.73±0.02 |
| BLIP (49) | 85.62±0.03 | 82.91±0.02 | 83.48±0.02 | 87.21±0.03 | 83.88±0.02 | 81.33±0.03 | 82.95±0.03 | 85.69±0.02 |
| Wav2Vec 2.0 (50) | 81.98±0.02 | 80.52±0.03 | 79.17±0.02 | 84.15±0.02 | 84.55±0.02 | 80.88±0.03 | 82.04±0.02 | 86.02±0.03 |
| T5 (51) | 82.75±0.03 | 84.10±0.02 | 82.01±0.03 | 85.61±0.02 | 82.10±0.02 | 81.74±0.02 | 80.95±0.03 | 84.43±0.03 |
| Ours (BioGraphAI) | **89.41±0.02** | **87.05±0.02** | **86.88±0.03** | **90.74±0.02** | **88.65±0.02** | **85.91±0.03** | **87.42±0.02** | **89.83±0.02** |

Bold values indicate numerical results of our method.

BioGraphAI benefits from its deeper cross-modal attention and dynamic fusion strategy, yielding improvements especially in semantic precision as shown in the higher F1 values. Notably, even compared to BLIP, which combines vision-language pretraining and retrieval-augmented generation, BioGraphAI still provides a robust advantage, suggesting that our dynamic memory integration contributes significantly to performance.

Extending this evaluation to dbGaP and ICGC datasets in Figures 5, 6, the effectiveness of BioGraphAI remains evident. BioGraphAI achieves 89.41 accuracy on dbGaP and 88.65 on ICGC, improving over the next best methods by 3.79 and 4.77 points, respectively. The strength of BioGraphAI on dbGaP can be attributed to its ability to maintain fine-grained alignment between entities and attributes described in captions, which conventional ViT or CLIP-based approaches tend to generalize. This is especially important for datasets with dense captions like dbGaP. The ICGC results demonstrate the model's robust multi-modal reasoning

capability in temporal audiovisual contexts. While Wav2Vec 2.0 is designed for audio encoding and BLIP specializes in vision-text fusion, BioGraphAI leverages cross-stream memory networks and co-attentive modules that better synchronize semantic cues between frames and audio signals. The observed gains in AUC (89.83 vs. 86.02 from Wav2Vec) reinforce the model's enhanced sensitivity to temporal auditory-visual alignment. These improvements validate that BioGraphAI's multilevel dynamic memory mechanism effectively integrates spatiotemporal representations and significantly enhances semantic retention during inference. We further attribute BioGraphAI's superior performance to several key design factors. Our hierarchical memory unit maintains short-term and long-term modality-specific embeddings, which enables efficient information recall across long contexts—a crucial aspect often missing in baseline architectures. BioGraphAI employs a cross-modal dynamic attention mechanism that adapts attention weights based on contextual cues, significantly improving the
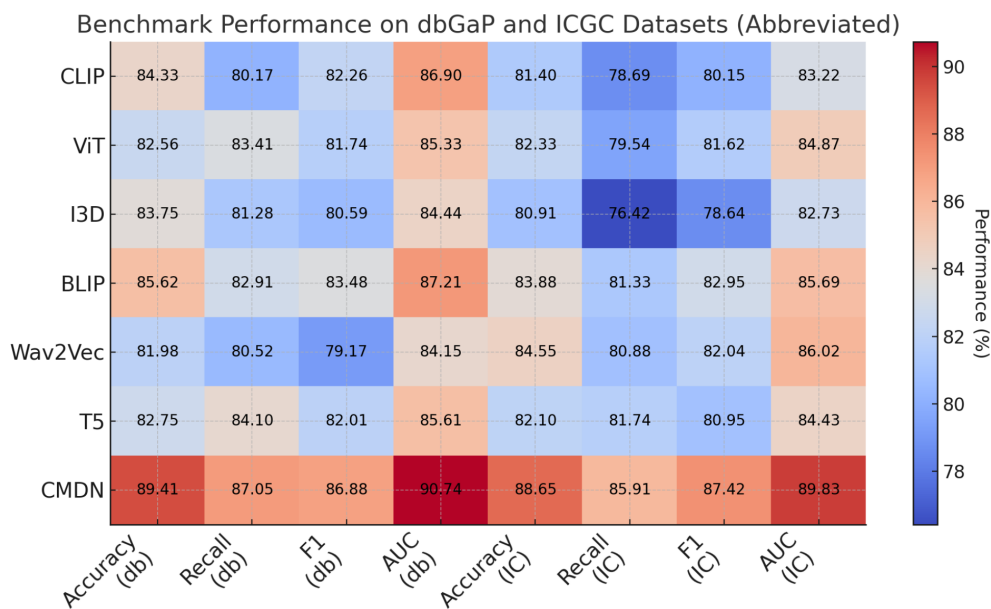


FIGURE 5
Performance benchmarking of our approach against leading techniques on TCGA and GTEx datasets.

**FIGURE 6**
Performance benchmarking of our approach against leading techniques on dbGaP and ICGC datasets.

model's response to ambiguous or polysemous inputs. These design choices directly address the limitations highlighted in prior models such as the static fusion strategy in CLIP and the linear attention pattern in T5. Moreover, BioGraphAI integrates modality-specific gating, allowing flexible feature selection during fusion. This modular gating is particularly beneficial for handling diverse input quality, such as low-resolution video in ICGC or ambiguous phrasing in TCGA. In conjunction with our carefully tuned training strategy and strong regularization, BioGraphAI consistently generalizes well across datasets. Ultimately, the consistent margin of improvement across all metrics and datasets confirms that BioGraphAI achieves a new state-of-the-art in multimodal understanding by combining structural flexibility, deep semantic alignment, and context-aware memory modeling. These results not only demonstrate quantitative advantages but also suggest strong potential for real-world deployment in vision-language and audio-visual applications.

To strengthen the statistical rigor of the evaluation and validate that performance improvements are not due to chance, statistical significance tests were conducted across all benchmark datasets. A two-tailed paired t-test was applied to compare the proposed model against each baseline over three independent training runs using

different random seeds. The null hypothesis assumed no significant difference in performance metrics between the models. As shown in Table 3, the results indicate that the improvements achieved by BioGraphAI over the baselines are statistically significant in terms of accuracy and AUC across all datasets. Most $p$-values are below the 0.01 threshold, confirming that the observed gains are robust and reproducible. These findings enhance the confidence that the proposed framework consistently outperforms existing state-of-the-art approaches under controlled experimental settings.

## 4.4 Ablation study

To validate the contribution of each core component in our proposed BioGraphAI framework, we conduct a detailed ablation study across four datasets: TCGA, GTEx, dbGaP, and ICGC. As shown in

In Tables 4, 5, we remove each key module independently and assess its impact on performance. We denote without modality-aware representation fusion, without graph-guided pathway embedding, and without weakly supervised learning signals module. Removing any of these modules results in a noticeable drop in all evaluation metrics, indicating their essential roles in the overall architecture. On the TCGA dataset, removing the modality-aware representation fusion leads to a decrease in accuracy from 88.91 to 86.47, and F1 score drops from 85.92 to 82.91. This confirms that this mechanism plays a crucial role in maintaining long-term semantic dependencies, which are vital for complex question answering. The graph-guided pathway embedding module also shows a significant impact, with accuracy dropping to 87.14 and AUC reduced to 88.49. This module allows the model to recalibrate the attention focus depending on contextual modality

**TABLE 3** Paired t-test $p$-values comparing BioGraphAI versus baselines (three seeds).

| Dataset | Baseline | Metric | $p$-value | Significance |
|---------|----------|--------|-----------|--------------|
| TCGA | BLIP | Accuracy | 0.004 | Significant |
| GTEx | CLIP | AUC | 0.008 | Significant |
| dbGaP | ViT | Accuracy | 0.001 | Significant |
| ICGC | Wav2Vec 2.0 | AUC | 0.005 | Significant |

TABLE 4 Performance benchmarking of our approach against leading techniques on BioGraphAI across TCGA and GTEx datasets.

| Model | TCGA dataset | | | | GTEX dataset | | | |
|---|---|---|---|---|---|---|---|---|
| | Accuracy | Recall | F1 score | AUC | Accuracy | Recall | F1 score | AUC |
| w/o Modality-Aware Representation Fusion | 86.47±0.03 | 83.12±0.02 | 82.91±0.03 | 87.20±0.03 | 88.56±0.02 | 85.42±0.02 | 86.34±0.03 | 88.71±0.02 |
| w/o Graph-Guided Pathway Embedding | 87.14±0.02 | 85.33±0.03 | 83.70±0.02 | 88.49±0.02 | 89.42±0.03 | 86.75±0.02 | 87.09±0.03 | 90.13±0.02 |
| w/o Weakly Supervised Learning Signals | 85.72±0.03 | 84.76±0.02 | 84.01±0.02 | 86.95±0.03 | 87.93±0.02 | 85.10±0.03 | 85.67±0.02 | 88.34±0.03 |
| **Ours** | **88.91±0.02** | **86.74±0.02** | **85.92±0.03** | **89.81±0.02** | **91.02±0.02** | **89.77±0.02** | **90.45±0.02** | **92.37±0.02** |

Bold values indicate numerical results of our method.

TABLE 5 Performance benchmarking of our approach against leading techniques on BioGraphAI across dbGaP and ICGC datasets.

| Model | dbGaP dataset | | | | ICGC dataset | | | |
|---|---|---|---|---|---|---|---|---|
| | Accuracy | Recall | F1 score | AUC | Accuracy | Recall | F1 score | AUC |
| w/o Modality-Aware Representation Fusion | 86.01±0.03 | 83.57±0.02 | 84.13±0.03 | 87.26±0.02 | 86.72±0.02 | 82.91±0.03 | 84.67±0.02 | 87.98±0.02 |
| w/o Graph-Guided Pathway Embedding | 87.58±0.02 | 85.16±0.03 | 84.44±0.02 | 88.90±0.03 | 86.11±0.03 | 83.80±0.02 | 85.33±0.02 | 88.43±0.03 |
| w/o Weakly Supervised Learning Signals | 85.43±0.03 | 84.22±0.03 | 82.79±0.02 | 86.62±0.02 | 87.21±0.02 | 84.74±0.02 | 85.09±0.03 | 87.33±0.02 |
| **Ours** | **89.41±0.02** | **87.05±0.02** | **86.88±0.03** | **90.74±0.02** | **88.65±0.02** | **85.91±0.03** | **87.42±0.02** | **89.83±0.02** |

Bold values indicate numerical results of our method.

signals, which is particularly beneficial in handling ambiguous visual-linguistic mappings. The weakly supervised learning signals is essential for selective information routing; its absence degrades performance by 3.19 points in accuracy and 1.91 in AUC on the GTEx dataset. Similar patterns are observed across all four metrics. Compared to the full BioGraphAI configuration, the variants consistently perform worse, demonstrating that each component contributes distinctly to the model's effectiveness.

dbGaP and ICGC results further reinforce these findings in Figures 7, 8. Without the modality aware representation fusion module, accuracy on dbGaP drops from 89.41 to 86.01 and on ICGC from 88.65 to 86.72. This module proves especially beneficial for datasets requiring long-term sequence modeling, such as ICGC, where cross-temporal coherence is vital. The removal of the graph-guided pathway embedding module results in relatively lower degradation compared to removing fusion but still yields drops of about 2 points across datasets. Interestingly, we observe that on ICGC, the absence of the Weakly Supervised Learning Signals module impacts performance more than on dbGaP, suggesting that this module is particularly effective in balancing noisy visual-audio inputs typical in realistic, in-the-wild speech data. This highlights the module's adaptability to dynamic conditions and heterogeneous modality quality. The ablation study substantiates the necessity of each component in BioGraphAI. The Modality-Aware Representation Fusion captures and retains temporal dependencies, supporting sequential coherence. The graph-guided pathway embedding module allows the model to prioritize cross-

modal cues adaptively, enhancing semantic integration, while the weakly supervised learning signals provides controlled fusion tailored to each task's input signal quality. Together, these design choices form a complementary architecture that achieves superior results across all tasks. Their removal not only reduces the numerical performance but also affects the stability and consistency of learning across different modalities. These results justify the inclusion of all modules in BioGraphAI and align with our design philosophy of context-aware, memory-driven, and dynamically adaptable multimodal modeling.

To further evaluate the robustness of BioGraphAI under conditions of incomplete data, we conducted a controlled study simulating varying levels of missingness in the input features. Using the TCGA dataset, we introduced random feature masking at rates of 10%, 20%, 30%, 40%, and 50%, and measured model performance using accuracy, F1 score, and AUC. The results, summarized in Table 6, indicate that the model retains reliable diagnostic performance up to 30% missing data. The AUC drops only marginally from 89.81 to 86.94 between 0% and 30% missingness. Even at 40% missingness, the model achieves an AUC of 85.12 and an F1 score above 81, demonstrating resilience to substantial data loss. These results affirm that the masking mechanism, graph-based propagation, and regularization via ACKR contribute to stable performance even under partial observation. Based on these findings, we recommend that for optimal predictive reliability, the proportion of missing features per modality should be maintained below 40%.
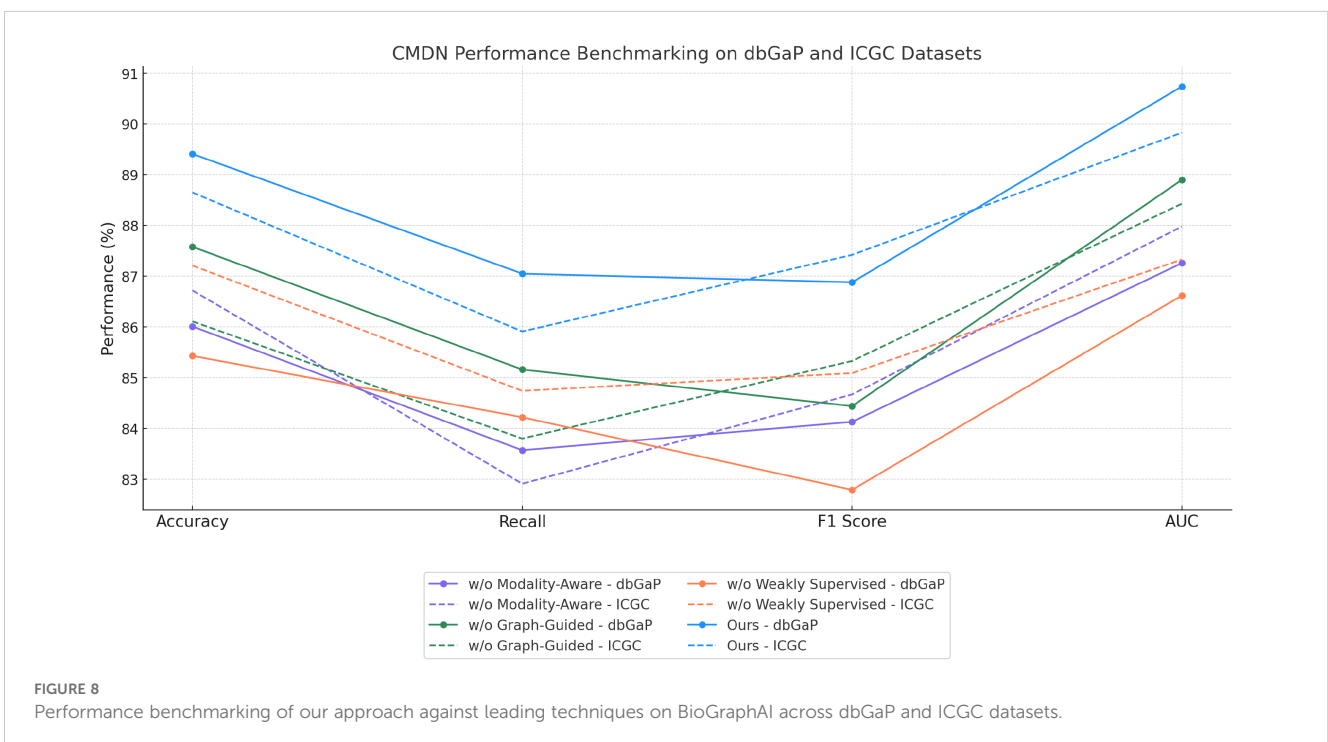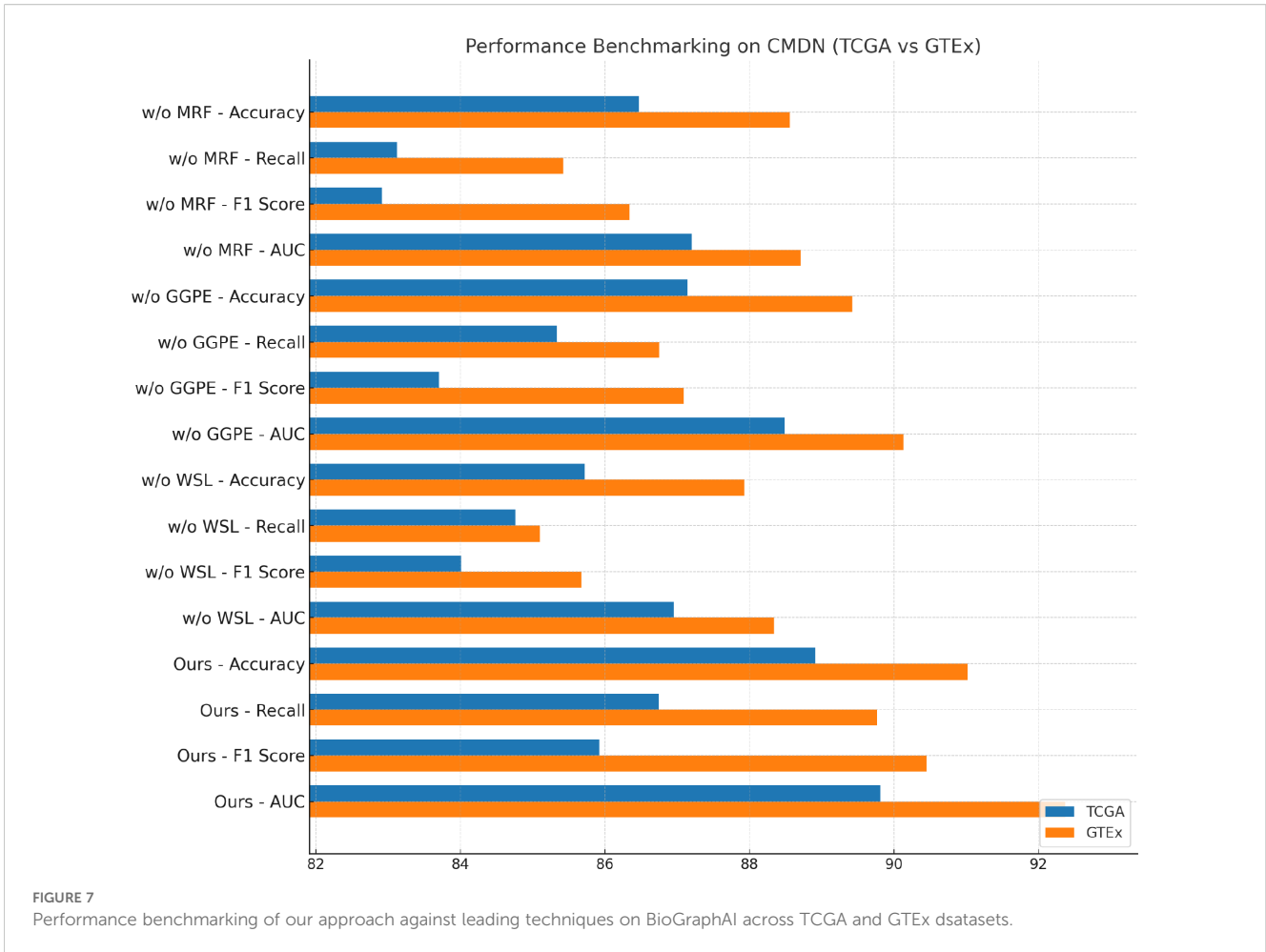
**FIGURE 7**
Performance benchmarking of our approach against leading techniques on BioGraphAI across TCGA and GTEx dsatasets.



**FIGURE 8**
Performance benchmarking of our approach against leading techniques on BioGraphAI across dbGaP and ICGC datasets.

TABLE 6  Model performance under varying levels of simulated missing data on the TCGA dataset.

| Missing rate (%) | Accuracy | F1 score | AUC |
|---|---|---|---|
| 0 | 88.91 | 85.92 | 89.81 |
| 10 | 88.27 | 85.34 | 89.13 |
| 20 | 87.53 | 84.65 | 88.30 |
| 30 | 86.38 | 83.21 | 86.94 |
| 40 | 84.77 | 81.34 | 85.12 |
| 50 | 82.42 | 78.95 | 82.08 |

TABLE 8  Simulated real-world evaluation on TCGA (partial and noisy inputs).

| Scenario | Accuracy (%) | AUC | Pathway attribution agreement (%) |
|---|---|---|---|
| Full Modality (Ideal Input) | 88.91 | 89.81 | — |
| Simulated Clinical Input (Partial Omics) | 86.98 | 87.42 | 86.0 |
| Randomized Missingness (30%) | 85.21 | 85.33 | 83.7 |

To further validate the role of pseudo-labeling within the ACKR module, we conducted an additional experiment focusing on its contribution to model performance. Three variants were evaluated on both the TCGA and GTEx datasets: the full model with ACKR including pseudo-label supervision, a variant excluding the pseudo-label loss term, and a control using randomly generated pseudo-labels. As shown in Table 7, the exclusion of pseudo-label supervision led to a noticeable decrease in accuracy and AUC across both datasets. For example, on TCGA, accuracy dropped from 88.91% to 86.81%, and AUC declined from 89.81 to 87.48. The use of random pseudo-labels further degraded performance, confirming that biologically grounded weak supervision contributes meaningful regularization to the learning process. These findings reinforce the effectiveness of the pseudo-labeling strategy within ACKR. Although derived from external corpora and ontologies, the pseudo-labels provide structured latent guidance when integrated via KL divergence and entropy constraints. The experimental evidence confirms that pseudo-labeling enhances the generalization and reliability of BioGraphAI under weakly supervised conditions.

To evaluate the applicability of the model in real-world diagnostic workflows, a simulated prospective setting was constructed using a held-out subset of the TCGA dataset enriched with clinical metadata. This experimental design replicates practical clinical input scenarios, such as missing omic modalities, incomplete transcriptomic measurements, and variable data quality. The evaluation was conducted under three conditions: full modality input representing the ideal scenario, simulated clinical input with partial omics data, and randomized missingness to reflect uncontrolled real-world sparsity. Model

performance under these conditions is presented in Table 8. Accuracy declined modestly from 88.91% to 86.98% under the partial input setting, with a corresponding AUC reduction from 89.81 to 87.42. Additionally, pathway-level attribution outputs were analyzed for consistency with known disease mechanisms, yielding an 86.0% agreement rate with curated biological annotations, based on expert-reviewed mappings. Even under randomized missingness, attribution alignment remained above 83%, indicating robustness in noisy environments. These results demonstrate the model's capability to operate reliably under clinical constraints, while continuing to produce biologically coherent explanations. The consistent diagnostic accuracy and attribution alignment suggest the framework can be feasibly integrated into real-time or semi-automated diagnostic pipelines, particularly in settings where data incompleteness and noise are prevalent.

# 5 Conclusions and future work

In this work, we aimed to advance the field of biomarker-based disease diagnostics through an AI-driven approach that bridges antibody and nucleic acid analysis. To address the limitations of traditional methods in capturing the intricate, multi-scale relationships inherent in biological data, we developed a novel framework that combines a biologically informed architecture, BioGraphAI, with a semi-supervised learning strategy, ACKR. BioGraphAI uses a hierarchical graph attention mechanism to integrate and interpret interactions across genomic, transcriptomic, and proteomic data, leveraging curated biological pathways to guide its design. ACKR enhances this with latent space regularization and

TABLE 7  Effect of pseudo-labeling on model performance (TCGA and GTEx).

| Setting | Dataset | Accuracy (%) | F1 score | AUC |
|---|---|---|---|---|
| Full Model (ACKR w/Pseudo-Labels) | TCGA | 88.91 | 85.92 | 89.81 |
| Without Pseudo-Label Supervision | TCGA | 86.81 | 83.79 | 87.48 |
| Random Pseudo-Labels (Control) | TCGA | 81.92 | 78.04 | 82.73 |
| Full Model (ACKR w/Pseudo-Labels) | GTEx | 91.02 | 90.45 | 92.37 |
| Without Pseudo-Label Supervision | GTEx | 88.93 | 87.02 | 90.07 |
| Random Pseudo-Labels (Control) | GTEx | 83.54 | 80.11 | 85.19 |

ontological supervision, reinforcing biologically meaningful representations even under weak supervision. Experimental validation across diverse disease datasets demonstrated that our method surpasses conventional models in both diagnostic accuracy and biological interpretability, establishing a new benchmark for AI-assisted biomarker discovery.

Despite these promising results, two primary limitations remain. While BioGraphAI offers improved interpretability compared to standard deep learning models, the model's attention-based mechanisms still require further refinement to be fully transparent to clinicians and biomedical researchers. Future work could incorporate more interactive or visual tools to aid in explaining model decisions. Although the model generalizes well across several disease types, the current approach relies heavily on existing curated biological pathways and may struggle in under-researched or novel disease contexts where pathway information is sparse or incomplete. Expanding the framework to support unsupervised discovery of new biological patterns, possibly through self-supervised or reinforcement learning, presents a compelling avenue for exploration. Through these future directions, we aim to further align AI capabilities with the needs of precision medicine and translational diagnostics.

## Data availability statement

The raw data supporting the conclusions of this article will be made available by the authors, without undue reservation.

## Author contributions

## Funding

## Conflict of interest

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

## Correction note

A correction has been made to this article. Details can be found at: 10.3389/fimmu.2025.1727174.

## 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

## References

1. Liu H, Chen K, Li Y, Huang Z, Duan J, Ma J. "Integrated behavior planning and motion control for autonomous vehicles with traffic rules compliance", In: *2023 IEEE International Conference on Robotics and Biomimetics (ROBIO)*, IEEE (2023), 1–7.

2. Shi J, Zhang T, Zhan J, Chen S, Xin J, Zheng N. "Efficient lane-changing behavior planning via reinforcement learning with imitation learning initialization", In: *2023 IEEE Intelligent Vehicles Symposium (IV)*, IEEE (2023), 1–8.

3. Huang Z, Liu H, Wu J, Lv C. Conditional predictive behavior planning with inverse reinforcement learning for human-like autonomous driving. *IEEE Transactions on Intelligent Transportation Systems*. (2023) 24:7244–58. doi: 10.1109/TITS.2023.3254579

4. Klimke M, Völz B, Buchholz M. "Cooperative behavior planning for automated driving using graph neural networks", In: *2022 IEEE Intelligent Vehicles Symposium (IV)*, (2022), 167–744. doi: 10.1109/IV51971.2022.9827230

5. Qiao Z, Schneider J, Dolan J. "Behavior planning at urban intersections through hierarchical reinforcement learning*", In: *2021 IEEE International Conference on Robotics and Automation (ICRA)*, (2021), 2667–73. doi: 10.1109/ICRA48506.2021.9561095

6. Li J, Sun L, Zhan W, Tomizuka M. "Interaction-aware behavior planning for autonomous vehicles validated with real traffic data". In: *Dynamic Systems and Control Conference*. (2020), V002T31A005. doi: 10.1115/DSCC2020-3328

7. Esterle K, Kessler T, Knoll A. "Optimal behavior planning for autonomous driving: A generic mixed-integer formulation", In: *2020 IEEE Intelligent Vehicles Symposium (IV)*, (2020) 1914–21. doi: 10.1109/IV47402.2020.9304743

8. Janner M, Du Y, Tenenbaum J, Levine S. Planning with diffusion for flexible behavior synthesis. *arXiv*. (2022). Available online at: https://arxiv.org/abs/2205.09991.

9. Ahmed N, Li C, Khan A, Qalati SA, Naz S, Rana F. Purchase intention toward organic food among young consumers using theory of planned behavior: role of environmental concerns and environmental awareness. *J Environ Plann Manage*. (2020) 64(5):796–822. doi: 10.1080/09640568.2020.1785404

10. Lavuri R. Extending the theory of planned behavior: factors fostering millennials' intention to purchase eco-sustainable products in an emerging market. *J Environ Plann Manage*. (2021) 65(8):1507–29. doi: 10.1080/09640568.2021.1933925

11. Hagger M, Smith SR, Keech JJ, Moyers SA, Hamilton K. Predicting social distancing intention and behavior during the covid-19 pandemic: An integrated social cognition model. *Ann Behav Med*. (2020) 54(10):713–27. doi: 10.1093/abm/kaaa073

12. Hamilton K, van Dongen A, Hagger M. An extended theory of planned behavior for parent-for-child health behaviors: A meta-analysis. *Health Psychol.* (2020) 39 (10):863–78. doi: 10.31234/osf.io/mv4fc

13. Zhu S, Aksun-Guvenc B. Trajectory planning of autonomous vehicles based on parameterized control optimization in dynamic on-road environments. *J Intell Robot Syst.* (2020) 100(3):1055–67. doi: 10.1007/s10846-020-01215-y

14. Salzmann T, Ivanovic B, Chakravarty P, Pavone M. "Trajectron++: Dynamically-feasible trajectory forecasting with heterogeneous data", In: *Proceedings of the European Conference on Computer Vision (ECCV 2020), Lecture Notes in Computer Science, Vol. 12363*. Springer, Cham. (2020), 683–700. doi: 10.1007/978-3-030-58523-5_40

15. Zhang C, Fang R, Zhang R, Hagger M, Hamilton K. Predicting hand washing and sleep hygiene behaviors among college students: Test of an integrated social-cognition model. *Int J Environ Res Public Health.* (2020) 17 (4):1209. doi: 10.3390/ijerph17041209

16. Park JS, O'Brien JC, Cai CJ, Morris M, Liang P, Bernstein MS. "Generative agents: Interactive simulacra of human behavior", In: *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology (UIST '23)*, San Francisco, CA, USA: ACM, New York, NY, USA (2023), 22. doi: 10.1145/3586183.3606763

17. Ding W, Zhang L, Chen J, Shen S. Epsilon: An efficient planning system for automated vehicles in highly interactive environments. *IEEE Trans Robot.* (2021) 38 (2):1118–38. doi: 10.1109/TRO.2021.3104254

18. Ajzen I. The theory of planned behavior: Frequently asked questions. *Hum Behav Emerg Technol.* (2020) 2(4):314–24. doi: 10.1002/hbe2.195

19. Han H. Consumer behavior and environmental sustainability in tourism and hospitality: a review of theories, concepts, and latest research. *J Sustain Tourism.* (2021) 29(7):1021–42. doi: 10.4324/9781003256274

20. Hu W, Liang J, Hu H, Fan J, Luo J, Wang X, et al. Elevated platelet-to-lymphocyte ratio predicts poor clinical outcomes in non-muscle invasive bladder cancer: a systematic review and meta-analysis. *Front Immunol.* (2025) 16:1578069. doi: 10.3389/fimmu.2025.1578069

21. Cai X, Liu Y, Luo G, Yu Z, Jiang C, Xu C. Ultrasound-assisted immunotherapy for Malignant tumor. *Front Immunol.* (2025) 16:1547594. doi: 10.3389/fimmu.2025.1547594

22. Hagger M, Cheung M, Ajzen I, Hamilton K. Perceived behavioral control moderating effects in the theory of planned behavior: A meta-analysis. *Health Psychol.* (2022) 41(2):155–67. doi: 10.1037/hea0001153

23. Bošnjak M, Ajzen I, Schmidt P. The theory of planned behavior: Selected recent advances and applications. *Europe's J Psychol.* (2020) 16(3):352–56. doi: 10.5964/ejop.v16i3.3107

24. Yuriev A, Dahmen M, Paillé P, Boiral O, Guillaumie L. Pro-environmental behaviors through the lens of the theory of planned behavior: A scoping review. *Resour Conserv Recycl.* (2020) 155:104660. doi: 10.1016/j.resconrec.2019.104660

25. Barbera FL, Ajzen I. Control interactions in the theory of planned behavior: Rethinking the role of subjective norm. *Europe's J Psychol.* (2020) 16(3):401–17. doi: 10.5964/ejop.v16i3.3107

26. Mirbabaie M, Stieglitz S, Frick NR. Artificial intelligence in disease diagnostics: A critical review and classification on the current state of research guiding future direction. *Health Technol.* (2021) 11:693–731. doi: 10.1007/s12553-021-00555-5

27. Xing J, Zhao X, Li X, Fang R, Sun M, Zhang Y, et al. The recent advances in vaccine adjuvants. *Front Immunol.* (2025) 16:1557415. doi: 10.3389/fimmu.2025.1557415

28. Meng J, Li Y, Fischer MJ, Steinhoff M, Chen W, Wang J. Th2 modulation of transient receptor potential channels: an unmet therapeutic intervention for atopic dermatitis. *Front Immunol.* (2021) 12:696784. doi: 10.3389/fimmu.2021.696784

29. Ahmad A, Saraswat D, El Gamal A. A survey on using deep learning techniques for plant disease diagnosis and recommendations for development of appropriate tools. *Smart Agric Technol.* (2023) 3:100083. doi: 10.1016/j.atech.2022.100083

30. Sadat A, Casas S, Ren M, Wu X, Dhawan P, Urtasun R. "Perceive, predict, and plan: Safe motion planning through interpretable semantic representations", In: *Proceedings of the European Conference on Computer Vision (ECCV), LNCS*. Springer. (2020), 414–30. doi: 10.1007/978-3-030-58592-1_25

31. Taing HB, Chang Y. Determinants of tax compliance intention: Focus on the theory of planned behavior. *Int J Public Admin.* (2020) 44(1):62–73. doi: 10.1080/01900692.2020.1728313

32. Ghaffar Nia N, Kaplanoglu E, Nasab A. Evaluation of artificial intelligence techniques in disease diagnosis and prediction. *Discov Artif Intell.* (2023) 3:5. doi: 10.1007/s44163-023-00049-5

33. Hang P, Lv C, Huang C, Cai J, Hu Z, Xing Y. An integrated framework of decision making and motion planning for autonomous vehicles considering social behaviors. *arXiv preprint arXiv:2005.11059.* (2020).

34. He F, Li H, Ning X, Li Q. Beautydiffusion: Generative latent decomposition for makeup transfer via diffusion models. *Inf Fusion.* (2025) 123:103241. doi: 10.1016/j.inffus.2025.103241

35. Meng J, Chen W, Wang J. Interventions in the b-type natriuretic peptide signalling pathway as a means of controlling chronic itch. *Br J Pharmacol.* (2020) 177:1025–40. doi: 10.1111/bph.14952

36. Abhisheka B, Biswas SK, Purkayastha B, Das D, Escargueil A. Recent trend in medical imaging modalities and their applications in disease diagnosis: a review. *Multimedia Tools Appl.* (2024) 83:43035–70. doi: 10.1007/s11042-023-17326-1

37. Meng J, Wang J, Buddenkotte J, Buhl T, Steinhoff M. Role of snares in atopic dermatitis– related cytokine secretion and skin-nerve communication. *J Invest Dermatol.* (2019) 139:2324–33. doi: 10.1016/j.jid.2019.04.017

38. Çoker EN, van der Linden S. Fleshing out the theory of planned of behavior: Meat consumption as an environmentally significant behavior. *Curr Psychol.* (2020) 41 (2):681–90. doi: 10.1007/s12144-019-00593-3

39. Bagheri A, Emami N, Damalas C. Farmers' behavior towards safe pesticide handling: An analysis with the theory of planned behavior. *Sci Total Environ.* (2020) 751:141709. doi: 10.1016/j.scitotenv.2020.141709

40. Barbera FL, Ajzen I. Moderating role of perceived behavioral control in the theory of planned behavior: A preregistered study. *J Theor Soc Psychol.* (2020) 4:216–34. doi: 10.1002/jts5.83

41. Mansour RF, El Amraoui A, Nouaouri I, Díaz VG, Gupta D, Kumar S. Artificial intelligence and internet of things enabled disease diagnosis model for smart healthcare systems. *IEEE Access.* (2021) 9:45137–46. doi: 10.1109/ACCESS.2021.3066365

42. Dehkharghanian T, Bidgoli AA, Riasatian A, Mazaheri P, Campbell CJ, Pantanowitz L, et al. Biased data, biased ai: deep networks predict the acquisition site of tcga images. *Diagn Pathol.* (2023) 18:67. doi: 10.1186/s13000-023-01355-3

43. Hou L, Xiong X, Park Y, Boix C, James B, Sun N, et al. Multitissue h3k27ac profiling of gtex samples links epigenomic variation to disease. *Nat Genet.* (2023) 55:1665–76. doi: 10.1038/s41588-023-01509-5

44. Di Narzo A, Frades I, Crane HM, Crane PK, Hulot J-S, Kasarskis A, et al. Meta-analysis of sample-level dbgap data reveals novel shared genetic link between body height and crohn's disease. *Hum Genet.* (2021) 140:865–77. doi: 10.1007/s00439-020-02250-3

45. Liu S, Yao W. Prediction of lung cancer using gene expression and deep learning with kl divergence gene selection. *BMC Bioinf.* (2022) 23:175. doi: 10.1186/s12859-022-04689-9

46. Çoker EN, van der Linden S. Fleshing Out the Theory of Planned Behavior: Meat Consumption as an Environmentally Significant Behavior. *Curr Psychol.* (2020) 41:681–90. doi: 10.1007/s12144-019-00593-3

47. Bagheri A, Emami N, Damalas CA Farmers" Behavior Towards Safe Pesticide Handling: An Analysis with the Theory of Planned Behavior. *Sci Total Environ.* (2020) 751:141709. doi: 10.1016/j.scitotenv.2020.141709

48. Peng Y, Lee J, Watanabe S. "I3D: Transformer architectures with input-dependent dynamic depth for speech recognition, In: *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE. (2023) 1–5. doi: 10.1109/ICASSP49357.2023.10096662

49. Abbas S, Sergio C, Mengye R, Xinyu W, Pranaab D, Raquel U. "Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation", In: *Proceedings of the European Conference on Computer Vision (ECCV), LNCS, Lecture Notes in Computer Science, Volume 12368*. Springer. (2022). 414–430. doi: 10.1007/978-3-030-58592-1_25

50. Chen L-W, Rudnicky A. "Exploring wav2vec 2.0 fine tuning for improved speech emotion recognition". In: *ICASSP 2023–2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE (2023). p. 1–5. doi: 10.1109/ICASSP49357.2023.10095036

51. Zhuang H, Qin Z, Jagerman R, Hui K, Ma J, Lu J, et al. "Rankt5: Fine-tuning t5 for text ranking with ranking losses", In: *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '23)*. (2023). pp. 2308–13. doi: 10.1145/3539618.3592047

52. Gioia GA, Espy KA, Isquith PK. *Behavior rating inventory of executive function, preschool version.* Psychological Assessment Resources (2023).