# Artificial Intelligence in Biomedical Research: From Data Integration to Precision Medicine

Anonymous authors
Paper under double-blind review

## Abstract

This comprehensive review examines the transformative role of artificial intelligence in biomedical research, from foundational data integration to clinical applications. The paper explores how AI techniques facilitate multimodal data fusion across diverse biological data types, employing both traditional statistical methods and advanced deep learning architectures including variational autoencoders, graph neural networks, and transformer models. It evaluates AI applications in medical imaging, where convolutional neural networks have achieved remarkable diagnostic accuracy (up to 94% in COVID-19 detection) while enhancing segmentation and classification tasks across multiple imaging modalities. The review further investigates generative AI's impact on molecular design and drug discovery, highlighting transformer-based architectures like TransAntivirus that navigate vast chemical spaces to optimize therapeutic candidates. Finally, it examines AI-enabled precision medicine applications, including Clinical Decision Support Systems and federated learning approaches that balance analytical power with privacy preservation. Despite significant progress, implementation challenges persist, including data heterogeneity, model explainability, and ethical concerns regarding bias and privacy. The paper underscores the importance of developing interpretable AI systems that integrate seamlessly into clinical workflows while addressing regulatory, ethical, and economic considerations to realize the full potential of AI in advancing biomedical research and healthcare delivery.

## 1 AI-Driven Multimodal Data Integration in Biomedical Research

### 1.1 Techniques and Applications of Multimodal Data Fusion

The convergence of multiple high-dimensional biological data types has created unprecedented opportunities to comprehensively characterize complex biological systems, transcending the limitations of single-modality approaches. Biomedical data fusion strategies span a spectrum from classical statistical methods to advanced deep learning architectures, each offering distinct advantages for integrating heterogeneous data types. Traditional integration approaches include joint non-negative matrix factorization (jNMF), partial least square (PLS), canonical correlation analysis (CCA), and multiple kernel learning (MKL). Among these, sparse canonical correlation analysis (CCA) has demonstrated particular utility in identifying relationships between different data modal- ities while maintaining interpretability. For instance, sparse CCA analysis of laboratory results and radiomics features in COVID-19 patients revealed significant correlations (cor$(Xu1, Zv1) = 0.596$), linking elevated lactate dehydrogenase and acute phase reactants with specific radiomic signatures that indicated increased entropy and heterogeneity in lung imaging features. This integration revealed distinct clinical phenotypes from multimodal data, underscoring the value of statistical approaches in uncovering biologically meaningful patterns across data types.

Deep learning has emerged as a transformative paradigm for multimodal biomedical data integration, offering superior capacity to model complex nonlinear relationships both within and across modalities. Variational autoencoders (VAEs) stand at the forefront of generative modeling approaches, providing flexible designs that balance dimensionality reduction

1

with generative capabilities for applications such as data imputation, denoising, and joint embedding creation. Recent methodological innovations have enhanced VAE perfor- mance through various regularization strategies, including adversarial training, cycle-consistency, contrastive learning, and disentangled representation learning. Intermediate fusion strategies with joint representation learning have proven particularly effective for capturing the intricate regulatory dynamics of biological systems, as they enable gradual fusion of modal- ities at different depths within the architecture, more closely reflecting true biological relation- ships. These approaches excel at bridging the heterogeneity gap between diverse data types—such as imaging, genomics, and clinical records—by applying modality-specific network architectures before integration, thereby mimicking the multilevel reasoning process used in clinical diagnosis and prognosis. Graph neural networks (GNNs) have further expanded the integration toolkit, proving especially valuable for analyzing complex relation- ships in graph-structured biomedical data, while transformer architectures with attention mecha- nisms have revolutionized cross-modality interaction learning.

The practical applications of multimodal data integration span the entire spectrum of biomedical research and clinical practice. In oncology, the integration of radiological images with genomic profiles has enhanced prognostic prediction and patient stratification, while the fu- sion of pathology whole-slide images with genomic features using attention-based frameworks like MCAT and PORPOISE has improved biomarker discovery. For COVID-19 anal- ysis, cooperative learning combining clinical, laboratory, radiomics, and viral genome sequencing data achieved superior prediction performance (AUC = 0.87) compared to unimodal approaches. Beyond molecular omics data, modern integration frameworks increasingly incor- porate diverse modalities including histopathology slides, MRI and PET images, electronic health records, and biosignals from wearable devices. This broader integration en- ables more holistic views of biological processesand diseases, particularly infields such as oncology and neurodegenerative disorders, where linking molecular mechanisms to clinical manifestations supports precise disease subtyping and targeted therapy development. Despite these advances, significant challenges remain, including data heterogeneity, missing modalities, limited cohort sizes, privacy concerns, and model interpretation[1]. Future directions point toward founda- tion models inspired by large language models, which aim to unify multimodal inputs, incorporate medical knowledge, and support reasoning across diverse biomedical data types.

## 2 ADVANCED AI APPLICATIONS IN BIOMEDICAL ANALYSIS

### 2.1 DEEP LEARNING IN MEDICAL IMAGING AND DIAGNOSTICS

Deep learning has emerged as a transformative approach to medical image analysis, significantly enhancing diagnostic capabilities across multiple healthcare domains. Convolutional Neural Net- works (CNNs) have become the dominant architecture for medical image processing due to their exceptional ability to extract meaningful features from complex visual data. These networks excel particularly in 2D data analysis, demonstrating rapid learning capabilities and strong performance when provided with sufficient labeled data . The application of CNNs and specialized frameworks such as AlexNet, VGG, Inception, and ResNet has substantially im- proved the efficacy of human clinicians in medical image interpretation, enabling more accurate and efficient diagnoses .

Medical image classification has achieved remarkable accuracy through deep learning implementa- tions. For instance, on the Image CLEF benchmark containing a 31-class image dataset, evolved techniques demonstrated an average classification accuracy of 88%, representing an improvement of up to 11% compared to previous state-of-the-art methods using identical datasets . In the context of COVID-19 diagnosis, the DeTraC (Decompose, Transfer, Com- pose) deep convolutional neural network achieved an impressive 94% accuracy with a true positive rate of 100% in identifying positive COVID-19 cases from chest X-ray images . Beyond classification, deep learning has revolutionized object localization and segmentation in medical imaging. For anatomical structure localization, ConvNets equipped with spatial pyramid pooling can analyze sagittal, axial, and coronal slices from three-dimensional images. These ca- pabilities are essential for numerous computer-aided diagnostic applications spanning microscopy,

ultrasound, computed tomography, dermoscopy, magnetic resonance imaging, positron emission tomography, and X-ray modalities .

Deep learning has demonstrated remarkable versatility across diverse medical domains, particularly in oncology. In gastric cancer research, models like convolutional neural networks (CNNs) and artificial neural networks (ANNs) have achieved significant praise in their application, revolutionizing diagnostic approaches . For cancer detection and characterization, deep neural networks can extract patterns from histopathology images to infer genomic features without re- quiring tumor sequencing. Notable examples include Image2TMB, which predicts tumor mutation burden in lung cancer with accuracy comparable to large panel sequencing but with significantly less variance, and HE2RNA, which infers gene expression from histopathology images to deter- mine microsatellite instability status in colorectal cancer . These approaches enable genomic inference from routine histopathological slides, potentially eliminating the need for expensive molecular testing in some clinical scenarios .

The integration of multimodal data through deep learning frameworks has opened new frontiers in medical image analysis. Advanced models can combine clinical data (diagnostic test results, pathol- ogy reports), medical images (histopathology, computed tomography), and various omics data (ge- nomic, transcriptomic, and proteomic profiles) to generate comprehensive insights . Autoencoder architectures, comprising encoders that create low-dimensional representation vectors and decoders that reconstruct the original input, have proven particularly effective for mul- timodal learning. This architecture forces the model to encapsulate meaningful features from the input, demonstrating good generalizability and the unique ability to integrate different data modal- ities into a single "end-to-end optimized" model . Recent innovations like the Segment Anything Model (SAM) have further advanced medical image segmentation capabil- ities, demonstrating superior performance compared to other interactive methods in non-iterative prompting settings, with notably better overall performance even when not used in its optimal box- prompting mode .

Despite these promising advances, significant challenges remain in implementing deep learning for medical imaging diagnostics. Data scarcity presents a fundamental limitation, as medical imag- ing datasets are typically much smaller than those used for general computer vision problems . The "black box" nature of many deep learning models raises legal and ethical concerns, as healthcare professionals may be reluctant to rely on systems whose decision-making processes cannot be fully explained . Additionally, the computational re- sources required to train complex deep learning models can be prohibitively expensive . Addressing these challenges will be crucial for realizing the full potential of deep learning in medical imaging diagnostics and facilitating its integration into clinical workflows .

## 2.2 Generative AI for Molecular Design and Drug Discovery

While deep learning has revolutionized medical imaging, similar transformative advances are oc- curring in molecular design and drug discovery, where generative AI models are redefining tradi- tional approaches. The integration of artificial intelligence methodologies into the drug development pipeline has yielded meaningful enhancements in both efficiency and effectiveness, particularly with the emergence of large language models and generative AI technologies [10]. These approaches are particularly valuable for navigating the vast chemical space—estimated to contain $10_{60}$ drug-like molecules—where traditional virtual screening techniques cannot feasibly explore the entirety of potential therapeutic candidates .

Generative models employing transformer-based architectures have demonstrated remarkable capa- bilities in molecular design tasks. The TransAntivirus framework represents a novel data-driven self-supervised pretraining generative model capable of performing select-and-replace edits to or- ganic molecules, optimizing them for desired properties in antiviral drug candidate development . This approach leverages the International Union of Pure and Applied Chemistry (IUPAC) nomenclature rather than Simplified Molecular Input Line Entry System (SMILES), pro- viding human-readable and easily editable molecular representations that more closely align with chemists' knowledge-based design practices . Such representation choice is partic- ularly advantageous for analogue-based drug design, where modifications typically involve altering

3

Under review as a
conference paper at
ICAIS 2025

functional groups rather than individual atoms, enabling more intuitive manipulation of molecular structures. Evaluations of TransAntivirus have demonstrated superior perfor- mance compared to control models across multiple metrics including novelty, validity, uniqueness, and diversity of generated compounds.

Recent advances in structure-based generative design have further enhanced the field by incorpo- rating protein structural information into de novo molecule optimization. These approaches can be categorized based on whether they employ distribution learning or goal-directed optimization, and whether they explicitly or implicitly incorporate protein structure information into the genera- tive model . Structure-based approaches aim to maximize predicted on-target binding affinity of generated molecules, thereby increasing the likelihood of identifying viable drug candidates with desired therapeutic properties . This integration of structural information represents a significant advancement over earlier generative models that relied solely on small-molecule information for training and conditioning de novo molecule generators .

The future of molecular design and drug discovery increasingly points toward AI agent systems capable of skeptical learning and reasoning. These "AI scientists" combine human creativity and expertise with artificial intelligence's capacity to analyze large datasets, navigate hypothesis spaces, and execute repetitive tasks . Rather than replacing human researchers, these biomedical AI agents function collaboratively, integrating various AI models and biomedical tools with experimental platforms . Such systems employ large language models and generative models with structured memory capabilities for continual learning, while incorporating scientific knowledge, biological principles, and theories through specialized machine learning tools . The potential applications span from virtual cell simulation and programmable control of phenotypes to the design of cellular circuits and development of novel therapeutics . As these AI-driven approaches continue to mature, they promise to accelerate the tradi- tionally lengthy and resource-intensive drug discovery process, potentially addressing urgent health challenges such as emerging viral threats through rapid identification and optimization of candidate molecules .

## 3 FROM DISCOVERY TO APPLICATION

### 3.1 AI-ENABLED PRECISION MEDICINE AND PERSONALIZED HEALTHCARE

The evolution from broad-spectrum therapeutic approaches to precision medicine represents one of the most significant paradigm shifts in modern healthcare, with artificial intelligence serving as a critical enabler of this transition. The integration of AI methodologies with multimodal clini- cal data has accelerated the development of tailored treatment protocols that account for individ- ual patient variability in genes, environment, and lifestyle. The novel Drug Intelligence Science (DIS®) platform exemplifies this advancement by combining single-cell technology with AI and machine learning to gain high-resolution insights into cell biology, thereby facilitating the discov- ery of disease-relevant targets, high-quality drug candidates, and predictive biomarkers . This innovative approach provides unprecedented mechanistic understanding of human dis- eases and enables in-depth pharmacological profiling of drug candidates, significantly increasing the probability of success in drug development and therapeutic interventions .

In the clinical context, AI-powered Clinical Decision Support Systems (CDSSs) have demonstrated substantial value in personalizing patient care. These systems incorporate features such as risk level estimation, diagnosis recommendations, and tailored treatment suggestions, collectively con- tributing to more effective healthcare delivery . The integration of AI into clinical workflows has shown positive outcomes across various implementation models, with both electronic medical record-integrated and stand-alone CDSSs demonstrating benefits to healthcare providers . For instance, MilleDSS, an Italian CDSS, illustrates practical im- plementation with its four domains of general practitioner-software interaction covering clinical management, prescribing appropriateness, prevention strategies, and medical computerized steward- ship. Despite these advances, the economic valuation of personalized medicine approaches remains complex. Studies suggest that while many personalized medicine tests are rela- tively cost-effective, fewer have been found to be cost-saving, and many available or emerging tests still require economic evaluation . This economic dimension underscores the

4

Under review as a
conference paper at
ICAIS 2025

need for more evidence to inform decision-making and assessment of genomic priorities in health-care resource allocation .

The success of AI in precision medicine ultimately depends on clinician trust and adoption. Research confirms that both accuracy and understandability are crucial for fostering clinician trust in predictive CDSSs, with the degree of reliance on these systems within clinical workflows potentially influencing trust requirements . Addressing this challenge, recent advances in interpretable machine learning have enabled the development of models that not only provide pre- dictions but also identify features that drive these predictions, as demonstrated in studies predicting physiological and perceived stress in pregnant women . Furthermore, AI models have proven valuable in optimizing clinical trials, patient selection, appropriate dosing regimens, and biomarker identification—all critical components of the personalized medicine ecosystem. These applications hold promise for streamlining clinical investigations and improv- ing patient outcomes through more targeted therapeutic approaches.

The implementation of AI-enabled precision medicine also presents significant data privacy and ethical considerations. Federated learning has emerged as a promising solution to data-sharing challenges, allowing models to be trained across multiple institutions without centralizing sensi- tive patient data. This approach facilitates cross-institutional collaboration while preserving privacy, an essential consideration given that the development of AI algorithms typically requires extensive processing of big data in biobanks. Legal frameworks such as the EU's General Data Protection Regulation (GDPR) provide guidance on ensuring compliance with data protection requirements when handling various categories of health data. Ad- ditionally, the Medical Device Regulation (EU 2017/745) stipulates that clinical evidence must be provided for any software intended for medical purposes, necessitating rigorous epidemiological studies to validate the effectiveness of AI systems in clinical practice.

## 3.2 EXPLAINABILITY, TRUST, AND IMPLEMENTATION CHALLENGES

Despite the transformative potential of AI in precision medicine, significant challenges persist in its widespread implementation. The deployment of these technologies necessitates careful attention to acomplexdevelopmentprocessthatbalancesautomationwithhumanexpertise. Deeplearningmod-els, while powerful, are not general-purpose AI systems but specialized tools that extract patterns from inputs and compute probabilities of class labels, requiring both representative training data and an understanding of their inherent limitations.. This technical reality underscores the importance of fostering intuitive understanding of these models among domain experts in areas such as health, education, and agriculture to facilitate effective translation to practice.. The potential for bias represents a particularly pressing concern in biomedical AI applications, with multiple sources of bias including insufficient data, sampling bias, and the use of health-irrelevant features or race-adjusted algorithms. Addressing these challenges requires sophis- ticated debiasing approaches that can be broadly categorized as distributional (data augmentation, perturbation, reweighting, andfederatedlearning)oralgorithmic(unsupervisedrepresentationlearn- ing, adversarial learning, disentangled representation learning, and causality-based methods).

The rapidly growing scale and variety of biomedical data repositories further complicate implemen- tationbyraisingimportantprivacyconcernsthatconventionaldata-sharingframeworksinadequately address. Privacy-enhancing technologies (PETs) have emerged as promising so- lutions that safeguard sensitive data while enabling broader usage and analysis. These technologies facilitate data sharing across institutional boundaries through mechanisms that provide formal privacy guarantees, with statistical disclosure control (SDC) and differential pri- vacy (DP) representing two dominant frameworks that address the fundamental statistical problem of balancing disclosure risk against data utility . Despite their different formulations, both approaches share core statistical challenges in designing optimal release mecha- nisms that satisfy bounds on disclosure risk while maximizing analytical utility . Beyond technical considerations, successful AI integration requires holistic attention to eth- ical and regulatory requirements. A comprehensive perspective on applications, opportunities, and challenges from a programmatic viewpoint is essential for ethical and sustainable implementation of AI solutions in medical contexts . This multifaceted approach ensures that

AI-based algorithms enhance outcomes, quality, and efficiency while respecting the complex social, ethical, and regulatory landscape in which they operate .

## 4 CONCLUSION

Artificial intelligence has fundamentally transformed biomedical research by enabling unprecedented integration and analysis of heterogeneous data modalities. This review demonstrates AI's significant contributions across the biomedical spectrum—from multimodal data fusion using variational autoencoders and graph neural networks to enhanced diagnostic capabilities in medical imaging and accelerated therapeutic discovery through generative models. The practical implementations in precision medicine highlight AI's potential to personalize healthcare delivery while improving clinical decision-making.

However, critical limitations persist that require focused attention. The 'black box' nature of many deep learning models undermines clinician trust and regulatory compliance, while data scarcity and quality issues compromise model generalizability. Privacy concerns and potential algorithmic biases further complicate clinical implementation. Economic evaluations of AI-enabled precision medicine approaches remain inadequate, complicating healthcare resource allocation decisions.

Future research should prioritize developing inherently explainable AI architectures that maintain high performance while providing interpretable insights. Federated learning and privacy-enhancing technologies deserve further investigation to enable collaborative model training without compromising data security. Additionally, standardized frameworks for evaluating AI systems' economic impact and clinical utility are essential. Most importantly, the field must evolve toward collaborative human-AI systems that augment rather than replace clinical expertise, ensuring that technological advancement serves the fundamental goal of improving patient outcomes through more targeted, effective, and accessible healthcare interventions.

## REFERENCES

Ana R Bai~ao, Zhaoxiang Cai, Rebecca C Poulos, Phillip J Robinson, Roger R Reddel, Qing Zhong, Susana Vinga, and Emanuel Gonc¸alves. A technical review of multi-omics data integration methods: from classical statistical to deep generative approaches. *arXiv preprint arXiv:2501.17729*, 2025a.

Ana R Bai~ao, Zhaoxiang Cai, Rebecca C Poulos, Phillip J Robinson, Roger R Reddel, Qing Zhong, Susana Vinga, and Emanuel Gonc¸alves. A technical review of multi-omics data integration methods: from classical statistical to deep generative approaches. *arXiv preprint arXiv:2501.17729*, 2025b.

Hyunghoon Cho, David Froelicher, Natnatee Dokmai, Anupama Nandi, Shuvom Sadhuka, Matthew M Hong, and Bonnie Berger. Privacy-enhancing technologies in biomedical data science. *Annual review of biomedical data science*, 7(1):317–343, 2024.

Iacopo Cricelli, Ettore Marconi, and Francesco Lapi. Clinical decision support system (cdss) in primary care: from pragmatic use to the best approach to assess their benefit/risk profile in clinical practice. *Current Medical Research and Opinion*, 38(5):827–829, 2022.

Geoff Currie, K Elizabeth Hawk, Eric Rohren, Alanna Vial, and Ran Klein. Machine learning and deep learning in medical imaging: intelligent imaging. *Journal of medical imaging and radiation sciences*, 50(4):477–487, 2019.

Ahmet Gorkem Er, Daisy Yi Ding, Berrin Er, Mertcan Uzun, Mehmet Cakmak, Christoph Sadee, Gamze Durhan, Mustafa Nasuh Ozmen, Mine Durusu Tanriover, Arzu Topeli, et al. Multimodal data fusion using sparse canonical correlation analysis and cooperative learning: a covid-19 cohort study. *NPJ digital medicine*, 7(1):117, 2024.

Shanghua Gao, Ada Fang, Yepeng Huang, Valentina Giunchiglia, Ayush Noori, Jonathan Richard Schwarz, Yasha Ektefaie, Jovana Kondic, and Marinka Zitnik. Empowering biomedical discovery with ai agents. *Cell*, 187(22):6125–6151, 2024.

Klaudia Grechuta, Pedram Shokouh, Ahmad Alhussein, Dirk Müller-Wieland, Juliane Meyerhoff, JeremyGilbert, SnehaPurushotham, CatherineRolland, etal. Benefitsofclinicaldecisionsupport systems for the management of noncommunicable chronic diseases: targeted literature review. *Interactive Journal of Medical Research*, 13(1):e58036, 2024.

Hao Guan, Pew-Thian Yap, Andrea Bozoki, and Mingxia Liu. Federated learning for medical image analysis: A survey. *Pattern recognition*, 151:110424, 2024.

Arnulf Jentzen, Benno Kuckuck, and Philippe von Wurstemberger. Mathematical introduction to deep learning: Methods, implementations, and theory. *arXiv preprint arXiv:2310.20360*, 2023.

Junwei Liu, Xiaoping Cen, Chenxin Yi, Feng-ao Wang, Junxiang Ding, Jinyu Cheng, Qinhua Wu, Baowen Gai, Yiwen Zhou, Ruikun He, et al. Challenges in ai-driven biomedical multimodal data fusion and analysis. *Genomics, Proteomics & Bioinformatics*, 23(1):qzaf011, 2025.

Jiashun Mao, Jianmin Wang, Amir Zeb, Kwang-Hwi Cho, Haiyan Jin, Jongwan Kim, Onju Lee, Yunyun Wang, and Kyoung Tai No. Transformer-based molecular generative model for antiviral drug design. *Journal of chemical information and modeling*, 64(7):2733–2745, 2023.

Ada Ng, Boyang Wei, Jayalakshmi Jain, Erin A Ward, S Darius Tandon, Judith T Moskowitz, Sheila Krogh-Jespersen, Lauren S Wakschlag, and Nabil Alshurafa. Predicting the next-day perceived and physiological stress of pregnant women by using machine learning and explainability: algorithm development and validation. *JMIR mHealth and uHealth*, 10(8):e33850, 2022.

Ajita Paliwal, Smita Jain, Sachin Kumar, Pranay Wal, Madhusmruti Khandai, Prasanna Shama Khandige, Vandana Sadananda, Md Khalid Anwer, Monica Gulati, Tapan Behl, et al. Predictive modelling in pharmacokinetics: from in-silico simulations to personalized medicine. *Expert Opinion on Drug Metabolism & Toxicology*, 20(4):181–195, 2024.

Malhar Patel, Ittai Dayan, Elliot K Fishman, Mona Flores, Fiona J Gilbert, Michal Guindy, Eugene J Koay, Michael Rosenthal, Holger R Roth, and Marius G Linguraru. Accelerating artificial intelligence: How federated learning can protect privacy, facilitate collaboration, and improve outcomes. *Health informatics journal*, 29(4):14604582231207744, 2023.

Kathryn A Phillips, Julie Ann Sakowski, Julia Trosman, Michael P Douglas, Su-Ying Liang, and Peter Neumann. The economic value of personalized medicine tests: what we know and what we need to know. *Genetics in Medicine*, 16(3):251–257, 2014.

Simon JD Prince. *Understanding deep learning*. MIT press, 2023.

Jessica M Schwartz, Maureen George, Sarah Collins Rossetti, Patricia C Dykes, Simon R Minshall, Eugene Lucas, and Kenrick D Cato. Factors influencing clinician trust in predictive clinical decision support systems for in-hospital deterioration: qualitative descriptive study. *JMIR Human Factors*, 9(2):e33960, 2022.

Liang Schweizer. Drug intelligence science (dis®): Pioneering a high-resolution translational platform to enhance the probability of success for drug discovery and development. *Drug Discovery Today*, 28(11):103795, 2023.

Aleksandra Slavkovi´c and Jeremy Seeman. Statistical data privacy: A song of privacy and utility. *Annual Review of Statistics and Its Application*, 10(1):189–218, 2023.

Süren Richard Stahlschmidt, Benjamin Ulfenborg, and Jane Synnergren. Multimodal deep learning for biomedical data fusion: a review. *Briefings in bioinformatics*, 23(2):bbab569, 2022.

S Suganyadevi, V Seethalakshmi, and KrishnasamyBalasamy. Areview ondeep learningin medical image analysis. *International Journal of Multimedia Information Retrieval*, 11(1):19–38, 2022.

Morgan Thomas, Andreas Bender, and Chris de Graaf. Integrating structure-based approaches in generative molecular design. *Current Opinion in Structural Biology*, 79:102559, 2023.

Takis Vidalis. Artificial intelligence in biomedicine: a legal insight. *BioTech*, 10(3):15, 2021.

Under review as a
conference paper at
ICAIS 2025

Yifan Yang, Mingquan Lin, Han Zhao, Yifan Peng, Furong Huang, and Zhiyong Lu. A survey of recent methods for addressing ai fairness and bias in biomedicine. *Journal of Biomedical Informatics*, 154:104646, 2024.

Chaoran Yu and Ernest Johann Helwig. Artificial intelligence in gastric cancer: A translational narrative review. *Annals of translational medicine*, 9(3):269, 2021.

## Emerging Research Directions in AI-Empowered Biomedical Research

Based on the provided paper content, I have identified three promising research directions that represent significant opportunities for advancing AI applications in biomedical research. Each direction addresses critical gaps in current approaches while leveraging emerging technological capabilities.

### 1. Multimodal Foundation Models for Biomedical Research

What is the direction?

This research direction focuses on developing large-scale foundation models specifically designed for biomedical applications that can seamlessly integrate and reason across heterogeneous biomedical data types. Similar to how large language models revolutionized natural language processing, biomedical foundation models would serve as versatile platforms capable of processing and generating insights from genomics, proteomics, imaging, clinical records, and other modalities simultaneously.

What are the innovations and challenges?

Innovations: - Unified representation learning frameworks that bridge the heterogeneity gap between diverse biomedical data types - Pre-training strategies that incorporate domain-specific medical knowledge and relationships - Transfer learning capabilities that allow adaptation to multiple downstream tasks with minimal fine-tuning - Integration of attention mechanisms to identify cross-modality interactions and contextual relationships

Challenges: - Extreme data heterogeneity across modalities (e.g., structured EHR data vs. unstructured imaging data) - Computational requirements for training models on massive multimodal datasets - Missing modalities and incomplete data in real-world biomedical datasets - Difficulty in establishing ground truth for complex multimodal relationships - Ensuring interpretability of model predictions for clinical acceptance

Significance to existing research Current approaches to multimodal data integration in biomedicine typically employ specialized models for specific modality combinations or tasks. Foundation models would transform this landscape by offering a unified framework that can serve as a base for numerous downstream applications. This would accelerate research across multiple fields simultaneously, from disease subtyping to drug discovery, by enabling knowledge transfer across domains and reducing the need for modality-specific model development.

Suggested research steps

1. Develop novel architecture designs that efficiently handle the unique characteristics of different biomedical data types

2. Create specialized pre-training objectives that capture biological relationships across modalities

3. Establish benchmark datasets that span multiple biomedical modalities with well-defined ground truth

4. Design modular components that allow for missing modality handling through data imputation techniques

5. Implement interpretability mechanisms that provide biological insights alongside predictions

6. Validate model performance on specific downstream tasks like disease prediction and biomarker discovery

7. Build model distillation techniques to create lightweight versions for clinical deployment

### 2. AI Agents for Molecular Design and Drug Discovery

What is the direction?

8

This direction aims to develop autonomous AI agent systems that combine multiple AI capabilities (generative models, reinforcement learning, reasoning systems) to actively participate in the drug discovery process. These systems would move beyond passive prediction to actively propose, test, and refine hypotheses about molecular designs with minimal human intervention, functioning as "AI scientists" that collaborate with human researchers.

What are the innovations and challenges?

Innovations: - Integration of skeptical learning and reasoning capabilities that allow agents to question and validate their own hypotheses - Closed-loop systems connecting in silico predictions with automated experimental platforms - Structured memory architectures enabling continual learning from successes and failures - Incorporation of scientific knowledge, biological principles, and theoretical constraints in agent reasoning processes - Multi-objective optimization capabilities that balance efficacy, toxicity, synthesizability, and other drug properties

Challenges: - Creating reliable simulation environments that accurately predict real-world molecular behavior - Developing frameworks for effective human-AI collaboration in hypothesis generation - Ensuring reproducibility and reliability of agent-designed experiments - Bridging the gap between computational predictions and wet-lab validation - Managing the vast hypothesis space efficiently while avoiding common chemical pitfalls

Significance to existing research Current generative models for drug discovery typically focus on passive molecule generation or optimization for specific properties in isolation. AI agent systems represent a paradigm shift by actively driving the discovery process through hypothesis generation, testing, and refinement. This approach could dramatically accelerate the traditionally lengthy drug discovery pipeline by navigating the vast chemical space more efficiently than either humans or traditional computational methods alone.

Suggested research steps

1. Develop integrated frameworks that combine generative chemistry models with reasoning capabilities

2. Create simulation environments that accurately reflect pharmacological principles and constraints

3. Design protocols for agent-driven hypothesis generation with built-in validation mechanisms

4. Implement feedback loops connecting computational predictions with experimental validation

5. Establish metrics to evaluate the novelty, diversity, and biological relevance of agent-proposed molecules

6. Build specialized knowledge bases that encode domain expertise in chemistry and biology

7. Validate the system by targeting specific therapeutic areas with unmet medical needs

3. Privacy-Preserving Collaborative AI for Precision Medicine

What is the direction?

This research direction focuses on developing AI frameworks that enable cross-institutional collaboration for precision medicine while maintaining strict privacy guarantees. These frameworks would leverage privacy-enhancing technologies such as federated learning, differential privacy, and secure multi-party computation to allow model training across distributed datasets without centralizing sensitive patient information.

What are the innovations and challenges?

Innovations: - Federated learning architectures optimized for heterogeneous biomedical data types - Differential privacy mechanisms that provide formal privacy guarantees while preserving data utility-Securemulti-partycomputationprotocolsforcollaborativemodeltrainingacrossinstitutions - Privacy-preserving techniques for multimodal data integration in clinical settings - Distributed validation frameworks for model performance assessment

Challenges: - Performance degradation in privacy-preserving settings compared to centralized approaches - Statistical challenges in balancing disclosure risk against analytical utility - Regulatory compliance across different jurisdictions and healthcare systems - Non-uniform data distributions

Under review as a
conference paper at
ICAIS 2025

across institutions (data heterogeneity) - Computational overhead of secure protocols in resource-constrained clinical environments - Trust establishment among participating institutions

Significance to existing research While AI has demonstrated remarkable potential in precision medicine, its widespread adoption is limited by data siloing and privacy concerns. Privacy-preserving collaborative AI addresses this fundamental limitation by enabling model development on much larger and more diverse patient populations without compromising confidentiality. This approach could significantly accelerate the translation of AI advances into clinical practice by facilitating evidence generation at scale while respecting ethical and regulatory requirements.

Suggested research steps

1. Develop federated learning algorithms specifically optimized for heterogeneous clinical data

2. Design privacy budgeting frameworks that maximize utility while maintaining privacy guarantees

3. Create benchmark datasets and evaluation metrics for privacy-preserving biomedical AI

4. Implement efficient secure computation protocols suitable for resource-constrained environments

5. Establish governance frameworks for multi-institutional collaboration

6. Validate approaches on real-world use cases such as rare disease diagnosis or treatment response prediction

7. Develop tools to quantify and communicate privacy-utility tradeoffs to stakeholders

In summary, these three research directions represent complementary approaches to advancing AI in biomedical research. Multimodal foundation models offer the broadest potential impact by creating versatile platforms for diverse applications. AI agents for drug discovery could revolutionize therapeutic development but require more time to mature and validate. Privacy-preserving collaborative AI addresses immediate barriers to clinical implementation and could see faster adoption in practice. Together, these directions address the key challenges identified in the paper while leveraging emerging AI capabilities to transform biomedical research and healthcare delivery.

Under review as a
conference paper at
ICAIS 2025