

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) (Liu et al., 2025). Among these, sparse canonical correlation analysis (CCA) has demonstrated particular utility in identifying relationships between different data modalities while maintaining interpretability. For instance, sparse CCA analysis of laboratory results and radiomics features in COVID-19 patients revealed significant correlations ($\text{cor}(X_{u1}, Z_{v1}) = 0.596$), linking elevated lactate dehydrogenase and acute phase reactants with specific radiomic signatures that indicated increased entropy and heterogeneity in lung imaging features (Er et al., 2024). 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 (Stahlschmidt et al., 2022). Variational autoencoders (VAEs) stand at the forefront of generative modeling approaches, providing flexible designs that balance dimensionality reduction

054 with generative capabilities for applications such as data imputation, denoising, and joint embed-
055 ding creation (Baião et al., 2025a;b). Recent methodological innovations have enhanced VAE per-
056 formance through various regularization strategies, including adversarial training, cycle-consistency,
057 contrastive learning, and disentangled representation learning (Baião et al., 2025a;b). Intermediate
058 fusion strategies with joint representation learning have proven particularly effective for capturing
059 the intricate regulatory dynamics of biological systems, as they enable gradual fusion of modalities
060 at different depths within the architecture, more closely reflecting true biological relationships
061 (Stahlschmidt et al., 2022). These approaches excel at bridging the heterogeneity gap between di-
062 verse data types—such as imaging, genomics, and clinical records—by applying modality-specific
063 network architectures before integration, thereby mimicking the multilevel reasoning process used
064 in clinical diagnosis and prognosis (Stahlschmidt et al., 2022). Graph neural networks (GNNs)
065 have further expanded the integration toolkit, proving especially valuable for analyzing complex
066 relationships in graph-structured biomedical data, while transformer architectures with attention
067 mechanisms have revolutionized cross-modality interaction learning (Liu et al., 2025; Baião et al.,
068 2025a).

069 The practical applications of multimodal data integration span the entire spectrum of biomedical
070 research and clinical practice. In oncology, the integration of radiological images with genomic
071 profiles has enhanced prognostic prediction and patient stratification (Liu et al., 2025), while the fu-
072 sion of pathology whole-slide images with genomic features using attention-based frameworks like
073 MCAT and PORPOISE has improved biomarker discovery (Liu et al., 2025). For COVID-19 anal-
074 ysis, cooperative learning combining clinical, laboratory, radiomics, and viral genome sequencing
075 data achieved superior prediction performance (AUC = 0.87) compared to unimodal approaches (Er
076 et al., 2024). Beyond molecular omics data, modern integration frameworks increasingly incorporate
077 diverse modalities including histopathology slides, MRI and PET images, electronic health records,
078 and biosignals from wearable devices (Baião et al., 2025a;b). This broader integration enables
079 more holistic views of biological processes and diseases, particularly in fields such as oncology and
080 neurodegenerative disorders, where linking molecular mechanisms to clinical manifestations sup-
081 ports precise disease subtyping and targeted therapy development (Liu et al., 2025). Despite these
082 advances, significant challenges remain, including data heterogeneity, missing modalities, limited
083 cohort sizes, privacy concerns, and model interpretation[1]. Future directions point toward founda-
084 tion models inspired by large language models, which aim to unify multimodal inputs, incorporate
085 medical knowledge, and support reasoning across diverse biomedical data types (Liu et al., 2025).

086 2 ADVANCED AI APPLICATIONS IN BIOMEDICAL ANALYSIS

087 2.1 DEEP LEARNING IN MEDICAL IMAGING AND DIAGNOSTICS

088 Deep learning has emerged as a transformative approach to medical image analysis, significantly
089 enhancing diagnostic capabilities across multiple healthcare domains. Convolutional Neural Net-
090 works (CNNs) have become the dominant architecture for medical image processing due to their
091 exceptional ability to extract meaningful features from complex visual data. These networks excel
092 particularly in 2D data analysis, demonstrating rapid learning capabilities and strong performance
093 when provided with sufficient labeled data (Suganyadevi et al., 2022). The application of CNNs
094 and specialized frameworks such as AlexNet, VGG, Inception, and ResNet has substantially im-
095 proved the efficacy of human clinicians in medical image interpretation, enabling more accurate and
096 efficient diagnoses (Suganyadevi et al., 2022).

097 Medical image classification has achieved remarkable accuracy through deep learning implementa-
098 tions. For instance, on the Image CLEF benchmark containing a 31-class image dataset, evolved
099 techniques demonstrated an average classification accuracy of 88%, representing an improvement
100 of up to 11% compared to previous state-of-the-art methods using identical datasets (Suganyadevi
101 et al., 2022). In the context of COVID-19 diagnosis, the DeTraC (Decompose, Transfer, Com-
102 pose) deep convolutional neural network achieved an impressive 94% accuracy with a true positive
103 rate of 100% in identifying positive COVID-19 cases from chest X-ray images (Suganyadevi et al.,
104 2022). Beyond classification, deep learning has revolutionized object localization and segmentation
105 in medical imaging. For anatomical structure localization, ConvNets equipped with spatial pyramid
106 pooling can analyze sagittal, axial, and coronal slices from three-dimensional images. These ca-
107 pabilities are essential for numerous computer-aided diagnostic applications spanning microscopy,

108 ultrasound, computed tomography, dermoscopy, magnetic resonance imaging, positron emission
109 tomography, and X-ray modalities (Suganyadevi et al., 2022).

110
111 Deep learning has demonstrated remarkable versatility across diverse medical domains, particularly
112 in oncology. In gastric cancer research, models like convolutional neural networks (CNNs) and
113 artificial neural networks (ANNs) have achieved significant praise in their application, revolution-
114 izing diagnostic approaches (Yu & Helwig, 2021). For cancer detection and characterization, deep
115 neural networks can extract patterns from histopathology images to infer genomic features without
116 requiring tumor sequencing. Notable examples include Image2TMB, which predicts tumor muta-
117 tion burden in lung cancer with accuracy comparable to large panel sequencing but with significantly
118 less variance, and HE2RNA, which infers gene expression from histopathology images to determine
119 microsatellite instability status in colorectal cancer (Jentzen et al., 2023). These approaches en-
120 able genomic inference from routine histopathological slides, potentially eliminating the need for
expensive molecular testing in some clinical scenarios (Jentzen et al., 2023).

121
122 The integration of multimodal data through deep learning frameworks has opened new frontiers in
123 medical image analysis. Advanced models can combine clinical data (diagnostic test results, pathol-
124 ogy reports), medical images (histopathology, computed tomography), and various omics data (ge-
125 nomic, transcriptomic, and proteomic profiles) to generate comprehensive insights (Jentzen et al.,
126 2023). Autoencoder architectures, comprising encoders that create low-dimensional representation
127 vectors and decoders that reconstruct the original input, have proven particularly effective for mul-
128 timodal learning. This architecture forces the model to encapsulate meaningful features from the
129 input, demonstrating good generalizability and the unique ability to integrate different data modal-
130 ities into a single “end-to-end optimized” model (Jentzen et al., 2023). Recent innovations like
131 the Segment Anything Model (SAM) have further advanced medical image segmentation capabil-
132 ities, demonstrating superior performance compared to other interactive methods in non-iterative
133 prompting settings, with notably better overall performance even when not used in its optimal box-
prompting mode (Guan et al., 2024).

134
135 Despite these promising advances, significant challenges remain in implementing deep learning for
136 medical imaging diagnostics. Data scarcity presents a fundamental limitation, as medical imaging
137 datasets are typically much smaller than those used for general computer vision problems (Sug-
138 anyadevi et al., 2022). The “black box” nature of many deep learning models raises legal and ethical
139 concerns, as healthcare professionals may be reluctant to rely on systems whose decision-making
140 processes cannot be fully explained (Suganyadevi et al., 2022). Additionally, the computational re-
141 sources required to train complex deep learning models can be prohibitively expensive (Suganyadevi
142 et al., 2022). Addressing these challenges will be crucial for realizing the full potential of deep learn-
143 ing in medical imaging diagnostics and facilitating its integration into clinical workflows (Guan
144 et al., 2024).

145 2.2 GENERATIVE AI FOR MOLECULAR DESIGN AND DRUG DISCOVERY

146
147 While deep learning has revolutionized medical imaging, similar transformative advances are oc-
148 ccurring in molecular design and drug discovery, where generative AI models are redefining tradi-
149 tional approaches. The integration of artificial intelligence methodologies into the drug development
150 pipeline has yielded meaningful enhancements in both efficiency and effectiveness, particularly with
151 the emergence of large language models and generative AI technologies [10]. These approaches are
152 particularly valuable for navigating the vast chemical space—estimated to contain 10^{60} drug-like
153 molecules—where traditional virtual screening techniques cannot feasibly explore the entirety of
potential therapeutic candidates (Thomas et al., 2023).

154
155 Generative models employing transformer-based architectures have demonstrated remarkable capa-
156 bilities in molecular design tasks. The TransAntivirus framework represents a novel data-driven
157 self-supervised pretraining generative model capable of performing select-and-replace edits to or-
158 ganic molecules, optimizing them for desired properties in antiviral drug candidate development
159 (Mao et al., 2023). This approach leverages the International Union of Pure and Applied Chem-
160 istry (IUPAC) nomenclature rather than Simplified Molecular Input Line Entry System (SMILES),
161 providing human-readable and easily editable molecular representations that more closely align with
chemists’ knowledge-based design practices (Mao et al., 2023). Such representation choice is partic-
ularly advantageous for analogue-based drug design, where modifications typically involve altering

162 functional groups rather than individual atoms, enabling more intuitive manipulation of molecular
163 structures (Mao et al., 2023). Evaluations of TransAntivirus have demonstrated superior perfor-
164 mance compared to control models across multiple metrics including novelty, validity, uniqueness,
165 and diversity of generated compounds (Mao et al., 2023).

166 Recent advances in structure-based generative design have further enhanced the field by incorpo-
167 rating protein structural information into de novo molecule optimization. These approaches can
168 be categorized based on whether they employ distribution learning or goal-directed optimization,
169 and whether they explicitly or implicitly incorporate protein structure information into the genera-
170 tive model (Thomas et al., 2023). Structure-based approaches aim to maximize predicted on-target
171 binding affinity of generated molecules, thereby increasing the likelihood of identifying viable drug
172 candidates with desired therapeutic properties (Thomas et al., 2023). This integration of structural
173 information represents a significant advancement over earlier generative models that relied solely
174 on small-molecule information for training and conditioning de novo molecule generators (Thomas
175 et al., 2023).

176 The future of molecular design and drug discovery increasingly points toward AI agent systems
177 capable of skeptical learning and reasoning. These “AI scientists” combine human creativity and
178 expertise with artificial intelligence’s capacity to analyze large datasets, navigate hypothesis spaces,
179 and execute repetitive tasks (Gao et al., 2024). Rather than replacing human researchers, these
180 biomedical AI agents function collaboratively, integrating various AI models and biomedical tools
181 with experimental platforms (Gao et al., 2024). Such systems employ large language models and
182 generative models with structured memory capabilities for continual learning, while incorporating
183 scientific knowledge, biological principles, and theories through specialized machine learning tools
184 (Gao et al., 2024). The potential applications span from virtual cell simulation and programmable
185 control of phenotypes to the design of cellular circuits and development of novel therapeutics (Gao
186 et al., 2024). As these AI-driven approaches continue to mature, they promise to accelerate the
187 traditionally lengthy and resource-intensive drug discovery process, potentially addressing urgent
188 health challenges such as emerging viral threats through rapid identification and optimization of
189 candidate molecules (Mao et al., 2023).

190 3 FROM DISCOVERY TO APPLICATION

191 3.1 AI-ENABLED PRECISION MEDICINE AND PERSONALIZED HEALTHCARE

192 The evolution from broad-spectrum therapeutic approaches to precision medicine represents one of
193 the most significant paradigm shifts in modern healthcare, with artificial intelligence serving as a
194 critical enabler of this transition. The integration of AI methodologies with multimodal clinical data
195 has accelerated the development of tailored treatment protocols that account for individual patient
196 variability in genes, environment, and lifestyle. The novel Drug Intelligence Science (DIS®) plat-
197 form exemplifies this advancement by combining single-cell technology with AI and machine learn-
198 ing to gain high-resolution insights into cell biology, thereby facilitating the discovery of disease-
199 relevant targets, high-quality drug candidates, and predictive biomarkers (Schweizer, 2023). This
200 innovative approach provides unprecedented mechanistic understanding of human diseases and en-
201 ables in-depth pharmacological profiling of drug candidates, significantly increasing the probability
202 of success in drug development and therapeutic interventions (Schweizer, 2023).

203 In the clinical context, AI-powered Clinical Decision Support Systems (CDSSs) have demonstrated
204 substantial value in personalizing patient care. These systems incorporate features such as risk level
205 estimation, diagnosis recommendations, and tailored treatment suggestions, collectively contribut-
206 ing to more effective healthcare delivery (Grechuta et al., 2024). The integration of AI into clinical
207 workflows has shown positive outcomes across various implementation models, with both electronic
208 medical record-integrated and stand-alone CDSSs demonstrating benefits to healthcare providers
209 (Grechuta et al., 2024). For instance, MilleDSS, an Italian CDSS, illustrates practical implementa-
210 tion with its four domains of general practitioner-software interaction covering clinical management,
211 prescribing appropriateness, prevention strategies, and medical computerized stewardship (Cricelli
212 et al., 2022). Despite these advances, the economic valuation of personalized medicine approaches
213 remains complex. Studies suggest that while many personalized medicine tests are relatively cost-
214 effective, fewer have been found to be cost-saving, and many available or emerging tests still require
215 economic evaluation (Phillips et al., 2014). This economic dimension underscores the need for more

216 evidence to inform decision-making and assessment of genomic priorities in healthcare resource allocation (Phillips et al., 2014).

217
218
219 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 (Schwartz et al., 2022). Addressing this challenge, recent advances in interpretable machine learning have enabled the development of models that not only provide predictions but also identify features that drive these predictions, as demonstrated in studies predicting physiological and perceived stress in pregnant women (Ng et al., 2022). 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 (Paliwal et al., 2024). These applications hold promise for streamlining clinical investigations and improving patient outcomes through more targeted therapeutic approaches.

229 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 sensitive patient data (Patel et al., 2023). 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 (Vidalis, 2021). 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 (Vidalis, 2021). Additionally, 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 (Cricelli et al., 2022)

241 3.2 EXPLAINABILITY, TRUST, AND IMPLEMENTATION CHALLENGES

243 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 a complex development process that balances automation with human expertise. Deep learning models, 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 (Prince, 2023). 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 (Prince, 2023). 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 (Yang et al., 2024). Addressing these challenges requires sophisticated debiasing approaches that can be broadly categorized as distributional (data augmentation, perturbation, reweighting, and federated learning) or algorithmic (unsupervised representation learning, adversarial learning, disentangled representation learning, and causality-based methods) (Yang et al., 2024).

258 The rapidly growing scale and variety of biomedical data repositories further complicate implementation by raising important privacy concerns that conventional data-sharing frameworks inadequately address (Cho et al., 2024). Privacy-enhancing technologies (PETs) have emerged as promising solutions that safeguard sensitive data while enabling broader usage and analysis (Cho et al., 2024). These technologies facilitate data sharing across institutional boundaries through mechanisms that provide formal privacy guarantees, with statistical disclosure control (SDC) and differential privacy (DP) representing two dominant frameworks that address the fundamental statistical problem of balancing disclosure risk against data utility (Slavković & Seeman, 2023). Despite their different formulations, both approaches share core statistical challenges in designing optimal release mechanisms that satisfy bounds on disclosure risk while maximizing analytical utility (Slavković & Seeman, 2023). Beyond technical considerations, successful AI integration requires holistic attention to ethical 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 (Currie et al., 2019). 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 (Currie et al., 2019).

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

- 324 Klaudia Grechuta, Pedram Shokouh, Ahmad Alhussein, Dirk Müller-Wieland, Juliane Meyerhoff,
325 Jeremy Gilbert, Sneha Purushotham, Catherine Rolland, et al. Benefits of clinical decision support
326 systems for the management of noncommunicable chronic diseases: targeted literature review.
327 *Interactive Journal of Medical Research*, 13(1):e58036, 2024.
328
- 329 Hao Guan, Pew-Thian Yap, Andrea Bozoki, and Mingxia Liu. Federated learning for medical image
330 analysis: A survey. *Pattern recognition*, 151:110424, 2024.
- 331 Arnulf Jentzen, Benno Kuckuck, and Philippe von Wurstemberger. Mathematical introduction to
332 deep learning: Methods, implementations, and theory. *arXiv preprint arXiv:2310.20360*, 2023.
333
- 334 Junwei Liu, Xiaoping Cen, Chenxin Yi, Feng-ao Wang, Junxiang Ding, Jinyu Cheng, Qinhua Wu,
335 Baowen Gai, Yiwen Zhou, Ruikun He, et al. Challenges in ai-driven biomedical multimodal data
336 fusion and analysis. *Genomics, Proteomics & Bioinformatics*, 23(1):qzaf011, 2025.
- 337 Jiashun Mao, Jianmin Wang, Amir Zeb, Kwang-Hwi Cho, Haiyan Jin, Jongwan Kim, Onju Lee,
338 Yunyun Wang, and Kyoung Tai No. Transformer-based molecular generative model for antiviral
339 drug design. *Journal of chemical information and modeling*, 64(7):2733–2745, 2023.
340
- 341 Ada Ng, Boyang Wei, Jayalakshmi Jain, Erin A Ward, S Darius Tandon, Judith T Moskowitz, Sheila
342 Krogh-Jespersen, Lauren S Wakschlag, and Nabil Alshurafa. Predicting the next-day perceived
343 and physiological stress of pregnant women by using machine learning and explainability: algo-
344 rithm development and validation. *JMIR mHealth and uHealth*, 10(8):e33850, 2022.
- 345 Ajita Paliwal, Smita Jain, Sachin Kumar, Pranay Wal, Madhusmruti Khandai, Prasanna Shama
346 Khandige, Vandana Sadananda, Md Khalid Anwer, Monica Gulati, Tapan Behl, et al. Predic-
347 tive modelling in pharmacokinetics: from in-silico simulations to personalized medicine. *Expert
348 Opinion on Drug Metabolism & Toxicology*, 20(4):181–195, 2024.
349
- 350 Malhar Patel, Ittai Dayan, Elliot K Fishman, Mona Flores, Fiona J Gilbert, Michal Guindy, Eu-
351 gene J Koay, Michael Rosenthal, Holger R Roth, and Marius G Linguraru. Accelerating artificial
352 intelligence: How federated learning can protect privacy, facilitate collaboration, and improve
353 outcomes. *Health informatics journal*, 29(4):14604582231207744, 2023.
- 354 Kathryn A Phillips, Julie Ann Sakowski, Julia Trosman, Michael P Douglas, Su-Ying Liang, and
355 Peter Neumann. The economic value of personalized medicine tests: what we know and what we
356 need to know. *Genetics in Medicine*, 16(3):251–257, 2014.
357
- 358 Simon JD Prince. *Understanding deep learning*. MIT press, 2023.
- 359 Jessica M Schwartz, Maureen George, Sarah Collins Rossetti, Patricia C Dykes, Simon R Minshall,
360 Eugene Lucas, and Kenrick D Cato. Factors influencing clinician trust in predictive clinical
361 decision support systems for in-hospital deterioration: qualitative descriptive study. *JMIR Human
362 Factors*, 9(2):e33960, 2022.
363
- 364 Liang Schweizer. Drug intelligence science (dis®): Pioneering a high-resolution translational plat-
365 form to enhance the probability of success for drug discovery and development. *Drug Discovery
366 Today*, 28(11):103795, 2023.
- 367 Aleksandra Slavković and Jeremy Seeman. Statistical data privacy: A song of privacy and utility.
368 *Annual Review of Statistics and Its Application*, 10(1):189–218, 2023.
369
- 370 Sören Richard Stahlschmidt, Benjamin Ulfenborg, and Jane Synnergren. Multimodal deep learning
371 for biomedical data fusion: a review. *Briefings in bioinformatics*, 23(2):bbab569, 2022.
- 372 S Suganyadevi, V Seethalakshmi, and Krishnasamy Balasamy. A review on deep learning in medical
373 image analysis. *International Journal of Multimedia Information Retrieval*, 11(1):19–38, 2022.
374
- 375 Morgan Thomas, Andreas Bender, and Chris de Graaf. Integrating structure-based approaches in
376 generative molecular design. *Current Opinion in Structural Biology*, 79:102559, 2023.
377
- Takis Vidalis. Artificial intelligence in biomedicine: a legal insight. *BioTech*, 10(3):15, 2021.

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

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

385 APPENDIX

387 Emerging Research Directions in AI-Empowered Biomedical Research

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

392 1. Multimodal Foundation Models for Biomedical Research

393 What is the direction?

394 This research direction focuses on developing large-scale foundation models specifically designed
395 for biomedical applications that can seamlessly integrate and reason across heterogeneous biomed-
396 ical data types. Similar to how large language models revolutionized natural language processing,
397 biomedical foundation models would serve as versatile platforms capable of processing and gener-
398 ating insights from genomics, proteomics, imaging, clinical records, and other modalities simulta-
399 neously.

400 What are the innovations and challenges?

402 **Innovations:** - Unified representation learning frameworks that bridge the heterogeneity gap be-
403 tween diverse biomedical data types - Pre-training strategies that incorporate domain-specific med-
404 ical knowledge and relationships - Transfer learning capabilities that allow adaptation to multiple
405 downstream tasks with minimal fine-tuning - Integration of attention mechanisms to identify cross-
406 modality interactions and contextual relationships

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

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

418 Suggested research steps

- 419 1. Develop novel architecture designs that efficiently handle the unique characteristics of different
420 biomedical data types
- 421 2. Create specialized pre-training objectives that capture biological relationships across modalities
- 422 3. Establish benchmark datasets that span multiple biomedical modalities with well-defined ground
423 truth
- 424 4. Design modular components that allow for missing modality handling through data imputation
425 techniques
- 426 5. Implement interpretability mechanisms that provide biological insights alongside predictions
- 427 6. Validate model performance on specific downstream tasks like disease prediction and biomarker
428 discovery
- 429 7. Build model distillation techniques to create lightweight versions for clinical deployment
- 430
- 431

432 2. AI Agents for Molecular Design and Drug Discovery

433 What is the direction?

434 This direction aims to develop autonomous AI agent systems that combine multiple AI capabilities
435 (generative models, reinforcement learning, reasoning systems) to actively participate in the drug
436 discovery process. These systems would move beyond passive prediction to actively propose, test,
437 and refine hypotheses about molecular designs with minimal human intervention, functioning as “AI
438 scientists” that collaborate with human researchers.
439

440 What are the innovations and challenges?

441 **Innovations:** - Integration of skeptical learning and reasoning capabilities that allow agents to ques-
442 tion and validate their own hypotheses - Closed-loop systems connecting in silico predictions with
443 automated experimental platforms - Structured memory architectures enabling continual learning
444 from successes and failures - Incorporation of scientific knowledge, biological principles, and theo-
445 retical constraints in agent reasoning processes - Multi-objective optimization capabilities that bal-
446 ance efficacy, toxicity, synthesizability, and other drug properties
447

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

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

459 Suggested research steps

- 460 1. Develop integrated frameworks that combine generative chemistry models with reasoning capa-
461 bilities
- 462 2. Create simulation environments that accurately reflect pharmacological principles and constraints
- 463 3. Design protocols for agent-driven hypothesis generation with built-in validation mechanisms
- 464 4. Implement feedback loops connecting computational predictions with experimental validation
- 465 5. Establish metrics to evaluate the novelty, diversity, and biological relevance of agent-proposed
466 molecules
- 467 6. Build specialized knowledge bases that encode domain expertise in chemistry and biology
- 468 7. Validate the system by targeting specific therapeutic areas with unmet medical needs

472 3. Privacy-Preserving Collaborative AI for Precision Medicine

473 What is the direction?

474 This research direction focuses on developing AI frameworks that enable cross-institutional collabo-
475 ration for precision medicine while maintaining strict privacy guarantees. These frameworks would
476 leverage privacy-enhancing technologies such as federated learning, differential privacy, and secure
477 multi-party computation to allow model training across distributed datasets without centralizing sen-
478 sitive patient information.
479

480 What are the innovations and challenges?

481 **Innovations:** - Federated learning architectures optimized for heterogeneous biomedical data types
482 - Differential privacy mechanisms that provide formal privacy guarantees while preserving data
483 utility - Secure multi-party computation protocols for collaborative model training across institutions
484 - Privacy-preserving techniques for multimodal data integration in clinical settings - Distributed
485 validation frameworks for model performance assessment

486 **Challenges:** - Performance degradation in privacy-preserving settings compared to centralized ap-
487 proaches - Statistical challenges in balancing disclosure risk against analytical utility - Regulatory
488 compliance across different jurisdictions and healthcare systems - Non-uniform data distributions
489 across institutions (data heterogeneity) - Computational overhead of secure protocols in resource-
490 constrained clinical environments - Trust establishment among participating institutions

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

497 **Suggested research steps**

- 498 1. Develop federated learning algorithms specifically optimized for heterogeneous clinical data
- 499 2. Design privacy budgeting frameworks that maximize utility while maintaining privacy guarantees
- 500 3. Create benchmark datasets and evaluation metrics for privacy-preserving biomedical AI
- 501 4. Implement efficient secure computation protocols suitable for resource-constrained environments
- 502 5. Establish governance frameworks for multi-institutional collaboration
- 503 6. Validate approaches on real-world use cases such as rare disease diagnosis or treatment response
- 504 7. Develop tools to quantify and communicate privacy-utility tradeoffs to stakeholders
- 505 prediction
- 506
- 507
- 508
- 509

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

517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539