*Research Article*

# Enhancing Clinical Research with Synthetic Patient Data: Leveraging ChatGPT for Improved Diagnostic Models

Mahnoor Mughal[1]*, Mohammad Irshad[2], Nisar Khan[3], Dr. Muhammad Iftikhar Hanif[4], Dr. Maqbool Khan[5], Nisa Irshad[6], Noor-Ul-Huda Waseem[7], Shafin I Chaudhri[8], Safa N Chaudhri[9], Khumal Butt[10], Aneeqa Mahmood[11]

[1]Quaid-i-Azam University & Researcher Genai-training.com LLC NYC, New York, USA

[2]AI LLM Technology Architect, Accenture, NYC, New York, USA

[3]Technology Development Director, Genai-trainings.com LLC, New York City, New York, USA

[4]Newcastle University Medicine Malaysia (NUMed Malaysia), Johor, Malaysia

[5]Assistant Professor, Pak-Austria Fachhoschule - Institute of Applied Sciences and Technology, Mang, Haripur, KPK, Pakistan

[6]Student, High School West; Half Hollow Hills School District New York

[7]St. George's University School of Medicine, St. George's, Grenada

[8]Student, Rutgers Preparatory School, Somerset, New Jersey, USA

[9]Pre-Medical Students New York Institute of Technology, Northern Boulevard, Valentines Ln, Old Westbury, New York, USA

[10]Mphil in Literature, Writer, Editor University of Central Punjab, Pakistan

[11]Australian College for Applied Psychology (ACAP) University College, Autralia

## A B S T R A C T

Recent strides in machine learning (ML) and generative artificial intelligence (GenAI) are transforming clinical research, opening new possibilities for privacy-conscious, data-driven insights into diagnosis, treatment and patient care. In clinical research, data scarcity, privacy concerns and limited access to high-quality datasets often hinder innovation. GenAI, leveraging synthetic data generated by advanced models like ChatGPT, addresses these challenges by emulating real patient histories without compromising patient privacy. This study examines how synthetic data generation can enhance ML applications in healthcare, simulating diverse clinical scenarios, supporting drug discovery and enabling personalized medicine. Using generative models such as generative adversarial network (GANs) and variational autoencoder (VAEs), this study explores the potential to produce synthetic data that maintains the integrity and statistical relevance of real-world data while safeguarding patient confidentiality. We also address the limitations and ethical concerns of AI-generated data, particularly around accuracy and interpretability. Our findings suggest that integrating synthetic data into clinical research could redefine healthcare practices by enabling scalable, privacy-preserving and more equitable access to medical insights.

By bridging the gap between real and synthetic patient data, this approach holds promise for advancing precision medicine and supporting evidence-based healthcare, ultimately fostering a transformative era in clinical research.

## 1. Introduction

Recent advancements in machine learning (ML) have revolutionized many industries, including healthcare, by enabling the analysis of vast amounts of clinical data to extract meaningful insights and improve decision-making processes.

The integration of ML in clinical research offers several key advantages. Firstly, ML models can assist in patient risk stratification and personalized treatment recommendations by analysing individual patient characteristics, genetic profiles and environmental factors[1]. This capability supports the shift towards precision medicine, where interventions are tailored each patient's specific needs to optimize therapeutic outcomes[2].

Secondly, ML facilitates the discovery of novel biomarkers and disease mechanisms through high-dimensional data analysis. By uncovering hidden patterns in biological signals or medical images, researchers can identify potential targets for therapeutic intervention and diagnostic innovation[3].

Moreover, ML-powered predictive models enhance clinical trial design and patient recruitment strategies, leading to more efficient studies with improved statistical power and reduced costs[4]. These models predict patient responses to treatments, stratify participants based on likelihood of treatment success and optimize trial protocols to maximize outcomes. It can also support drug discovery by modeling patient responses to drugs and enhance personalized treatment plans by reflecting individual patient characteristics. This approach has been demonstrated in various applications, such as medical imaging and clinical trial simulations[5-7].

Machine learning algorithms, including deep learning and natural language processing (NLP), have demonstrated remarkable capabilities in analysing diverse healthcare datasets, such as electronic health records (EHRs), medical imaging, genomic data and wearable sensor data[8,9]. These algorithms can identify patterns, predict outcomes and uncover associations that may not be apparent through traditional statistical methods alone. For instance, Desai et al.2017[10] used deep learning to predict sepsis risk in Intensive Care Unit (ICU) patients, showing improved predictive accuracy compared to traditional methods. Similarly, NLP tools have streamlined literature review and meta-analysis by automating the extraction and summarization of relevant studies. This has been particularly useful in conducting systematic reviews[11] and outcomes of medical interventions, has increasingly leveraged ML techniques to address complex challenges and enhance evidence-based practice. Issues related to data quality, interpretability of ML models, regulatory compliance and ethical considerations must be carefully addressed to ensure the reliability and ethical integrity of research findings[12]. Ensuring the privacy and security of sensitive patient data remains a major concern. Regulations such as the Health Insurance Portability and Accountability Act (HIPAA) in the United States impose strict guidelines for handling healthcare data, which can complicate the use of deep learning and NLP in clinical research[13,14].

As an alternative, in clinical research, the potential applications of generative AI are vast and transformative.

Generative AI (GenAI), a broad category of artificial intelligence systems, focuses on generating new content or data. This includes creating text, images, audio, video and other forms of digital media based on learned patterns from existing data. GenAI leverages various techniques and models to produce outputs that mimic human creativity or innovation. GenAI represents a transformative approach to create new data and content that emulates existing data distributions. Unlike traditional AI models focused on classification or regression, GenAI seeks to understand and replicate underlying patterns to produce novel, yet realistic data. This technology has gained considerable attention due to its potential to revolutionize various fields, including clinical research[15]. GenAI holds substantial promise for advancing clinical research by addressing several critical challenges:

**1.1. Data Scarcity and Privacy:** The scarcity of high-quality, annotated medical data is a well-known challenge [20]. Generative AI can create synthetic medical data that maintains the statistical properties of real data without compromising patient privacy. This capability is crucial for training robust machine learning models and validating new algorithms without the need for extensive real-world datasets.

Simulation of Clinical Scenarios: Accurate simulation of clinical scenarios can provide valuable insights into disease progression and treatment responses[16]. Generative models can simulate diverse patient populations and clinical conditions, allowing researchers to explore various scenarios and predict outcomes under different treatment regimens[17].

Drug Discovery and Development: Generative AI accelerates drug discovery by proposing novel molecular structures and predicting their interactions with biological targets. This approach can significantly shorten the drug development timeline and enhance the efficiency of discovering new therapeutic agents[18].

Personalized Medicine: By generating synthetic patient profiles, Generative AI supports the development of personalized treatment strategies tailored to individual patient characteristics. This capability enables researchers to better understand how different patients might respond to various treatments and optimize therapeutic approaches[19].

In clinical research, where data availability, privacy and ethical concerns often pose significant barriers[20], GenAI offers innovative solutions. Techniques such as Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs) are particularly noteworthy. GANs consist of two competing neural networks a generator and a discriminator-that work together to create realistic synthetic data, while VAEs learn probabilistic models of data and can generate new samples from this learned distribution. These methodologies have shown promise in enhancing data generation, simulation and modelling tasks in clinical research[21,22]. GANs have been used to create synthetic medical images that aid in training diagnostic algorithms without exposing patient data[23]. Similarly, VAEs can model complex biological systems and generate synthetic patient records to support personalized medicine research. Moreover, GenAI enables the simulation of clinical scenarios,

understanding disease progression, evaluating treatment strategies and predicting patient outcomes. This ability is crucial for developing robust and adaptable therapeutic approaches.

This paper aims to explore how synthetic data can revolutionize clinical research, offering hope for a future where medical advancements are both groundbreaking and compassionate. By using generative AI to create realistic, privacy-preserving patient data, it seeks to empower healthcare professionals to make more accurate diagnoses, discover new treatments and provide personalized care, all while ensuring that patient trust remains at the heart of every innovation. Ultimately, it strives to create a healthcare system where every step forward brings us closer to a world of better, more equitable care for all.

## 2. Literature Review

Synthetic data in healthcare has emerged as a promising solution to address privacy concerns while facilitating data sharing for research and innovation[24]. Research has highlighted several advanced techniques and models for generating synthetic health data, each aimed at overcoming barriers related to privacy, data availability and ethical considerations.

High-dimensional synthetic data can help navigate these challenges and methods like the Gaussian Copula and Tabular Variational Autoencoder have been proposed to ensure privacy by anonymizing patient information[25]. The development of deep generative models has further expanded the potential for creating realistic synthetic health datasets, which preserve key characteristics of real data without disclosing sensitive information, thus supporting the development of predictive models and health IT platforms[26]. Novel algorithms, such as MIIC-SDG, generate synthetic data based on multivariate information frameworks, effectively balancing data quality with privacy concerns[27].

In practical applications, synthetic health data generated from administrative records plays a significant role in drug safety studies. For instance, ModOSIM produces data that more closely resembles real-world records compared to OSIM2, thereby enhancing methodological research and analyst training. Studies have focused on generating synthetic health data for longitudinal cohort studies, demonstrating the ability of synthetic datasets to reproduce real-world analysis results, particularly in nutrition research[28]. Evaluations of synthetic data in these contexts include assessing variable distributions, correlations and dependencies, with real-world analysis results being largely reproducible.

Furthermore, GCP tensor decomposition models have been developed to generate high-quality synthetic longitudinal health data, preserving patient privacy while maintaining data similarity to real data[29]. These models ensure both the utility and privacy of the original data.

Reviews of synthetic data generation methods in healthcare discuss the current status of the field, highlighting advancements, techniques and the effectiveness of synthetic data as an alternative to real data in research. They also analyse the challenges and opportunities associated with synthetic data in healthcare[30]. Synthetic data is increasingly used to address issues of data availability, privacy and bias propagation in medical applications. It supports all stages of model development, including clinical risk prediction, without the need for real data access, thus enhancing collaboration and project efficiency

[30]. The potential of synthetic data for open-access healthcare datasets is noted, offering a promising future for privacy protection and data sharing in healthcare research[31,32]. In particular, synthetic data in women's health provides a valuable alternative for research, efficiently addressing challenges in obtaining real-world data, especially for epidemiological and clinical problems[32].

## 3. Limitations of the study

The fallibility of generative accuracy in AI models, such as ChatGPT, has raised significant concerns, particularly in the context of clinical settings. A study revealed that 36% of ChatGPT-generated documents contained erroneous information, which underscores the need for development with updated AI models to improve accuracy and reduce the potential for misinformation[33]. Moreover, there is poor agreement on the quality of patient histories generated by AI, with reviewers often differing in their assessments. This discrepancy is further complicated by the challenges associated with patient multimedia components, including voices, images, animations and videos, which often suffer from unrealistic representations, stereotypes and technical issues.

Additionally, the accuracy and validity of synthetic data produced by AI models are often questioned, with limited assessment of ChatGPT's performance in clinical environments. Concerns about the quality and reliability of information provided by these models are compounded by issues of authorship, bias and the risk of generating inaccurate content, which could lead to misinformation or even research fraud. The lack of updated datasets for literature reviews also poses a significant risk to the accuracy of AI-generated content, with some models, like GPT-3.5, found to lack currency and omit important information.

In patient education, ChatGPT-generated materials are often less understandable and readable than traditional reference handouts, making them less ideal for this purpose. These materials tend to be written at higher grade levels with longer sentences, potentially making them difficult for patients to comprehend. Furthermore, ChatGPT is not specifically designed for medical information and the materials it produces require rigorous editing and fact-checking to avoid spreading misinformation, misdiagnosis or bias. This could disrupt healthcare relationships and raise significant ethical and privacy concerns.

The philosophical nature of AI communication and intelligence in companionship also invites scrutiny, particularly when considering the factual inaccuracies, omissions and inaccurate assumptions that may arise. ChatGPT's lack of flexibility in recommendations and its focus on specific topics can limit the usefulness of its responses in diverse situations. Consequently, the potential for misuse and misinformation by generative AI emphasizes the need for continued research, refinement and a strong commitment to ethical standards to ensure the reliability and safety of AI in healthcare and other critical fields.

Researchers face multiple challenges in applying Large Language Models (LLMs) like GPT-4 and LLaMA to complex mathematical tasks. These models often lack accuracy and reliability in high-stakes contexts, producing inconsistent solutions. Additionally, they lack effective decision-making frameworks, struggle to balance exploration with exploitation and operate in an infinite action space, complicating precise

output. LLMs also lack efficient feedback mechanisms and are difficult to integrate with structured decision-making algorithms like Monte Carlo Tree Search (MCTS). Designing scoring and reward systems for improvement is challenging, as is managing the high computational cost required to ensure accuracy and performance.

## 4. Methodology

### 4.1. Data Collection and Baseline Generation

An anonymized sample of real clinical data is collected to establish baseline characteristics, including demographics, medical histories and treatments. This data is used to define parameters for generating synthetic patient histories that reflect essential clinical variations.

### 4.2. Synthetic Data Generation Using ChatGPT

ChatGPT generates synthetic patient histories based on carefully designed prompts and constraints that align with baseline parameters. The model is fine-tuned to ensure data diversity, capturing variations across age, disease type and medical history, thus simulating a wide range of clinical scenarios.

### 4.3. Validation and Comparative Analysis

The synthetic data is evaluated through statistical analysis to ensure its distribution and correlation align with real data. Clinical experts review the data for plausibility. Diagnostic models are then trained using both real and synthetic data and performance metrics (F1 score, ROC-AUC) are used to compare model accuracy, sensitivity and specificity.

### 4.4. Study gap

One significant gap is the limited discussion on patient feedback and comprehension, which leaves uncertainty regarding how well patients understand and benefit from AI-generated materials. Additionally, there is a lack of comparison between AI-generated patient reports and traditional ones, which is crucial for determining the effectiveness of AI in clinical settings. Concerns about privacy, bias and accuracy of synthetic patient data are also prominent, with ethical considerations regarding the risks and benefits to patients' data needing further investigation. Moreover, the role of AI in clinical diagnosis and treatment decisions remains uncertain, partly due to the lack of rigorous evaluation methods assessing the educational effectiveness of AI-generated materials. Furthermore, there is a noticeable absence of analysis on the quality and reliability of the information provided by AI, as well as a lack of comparison with other educational tools or resources.

## 5. Results

The document provides a detailed exploration of synthetic data and its transformative potential in clinical research through a series of structured tables and graphs. **(Table 1)** highlights the barriers in traditional clinical research, such as data scarcity, privacy concerns and high costs, juxtaposed with solutions offered by Generative AI (GenAI). These solutions include synthetic dataset generation for enhanced data availability, privacy preservation through anonymization and simulation of diverse clinical scenarios, underscoring the ability of GenAI to address longstanding challenges in healthcare research. **(Table 2)** further elaborates on synthetic data generation techniques, comparing methods like GANs, VAEs and Gaussian Copula by detailing their advantages, such as the realism of GAN-generated data and the interpretability of VAEs, alongside limitations like computational intensity and potential lack of diversity. It also highlights use cases such as disease progression modelling and privacy-preserving electronic health records (EHRs).

**5.1. A synthetic dataset is presented in the Table:** Synthetic Data of Breast Cancer Patients, which showcases simulated profiles of patients aged 40-65 years, categorized by ethnicity, gender, race, medical history, tumour stage and treatment plans. This dataset demonstrates how synthetic data can emulate diverse patient attributes while ensuring privacy. The corresponding Graph of Synthetic Data Visualization for Breast Cancer Patients visually represents this dataset, illustrating trends like age distribution, ethnicity and treatment preferences, thereby validating the utility of synthetic data in research and decision-making.

The Graph of Performance Comparison of Diagnostic Models compares the efficacy of two models, A and B, in terms of accuracy, sensitivity and specificity. It highlights the enhanced performance of models trained with diverse and balanced synthetic datasets, showcasing the practical applications of synthetic data in AI-driven diagnostics. **(Table 4)**: Specialists Evaluation of Synthetic Data Plausibility records the insights of clinical experts from various specialties, providing credibility ratings (1-5) along with observations on synthetic data's strengths and limitations. While some experts commend the data's realism and utility in training AI models, others note its shortcomings in capturing rare conditions or complex clinical scenarios.

Lastly, the **(Table5)**: Ethical Risk vs. Benefit Analysis of Synthetic Data Use addresses privacy concerns, bias and potential misuse, balanced against benefits like improved data accessibility, reduced bias in AI models and accelerated clinical trials. Mitigation strategies, such as robust anonymization and clear regulatory frameworks, are also outlined, emphasizing the importance of ethical considerations in deploying synthetic data.

**Table 1:** Summary of Key Challenges in Clinical Research and GenAI Solutions.
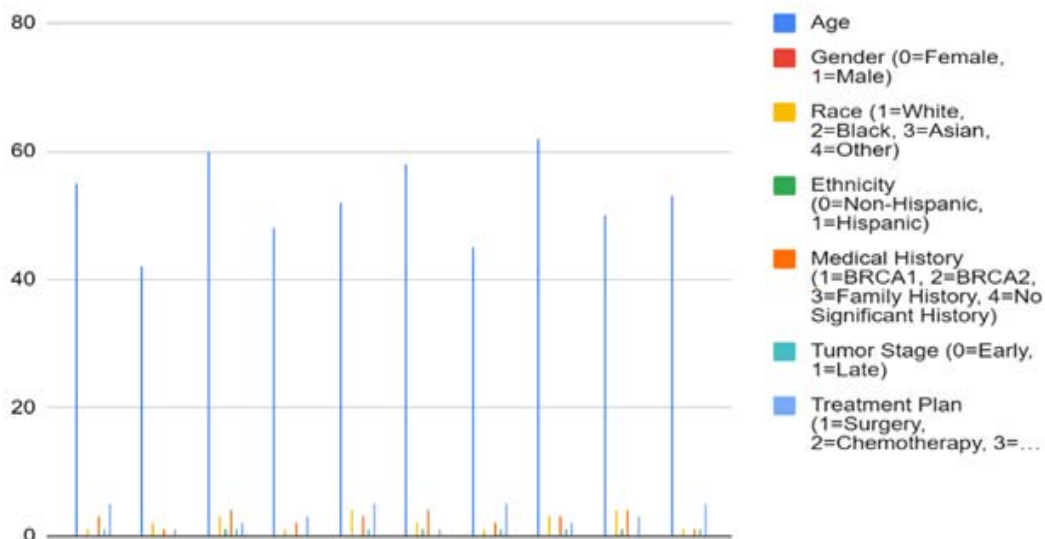
| Barriers in Traditional Clinical Research | Solutions Offered by Generative AI (GenAI) |
|---|---|
| Data Scarcity | Generates synthetic datasets that mimic real-world data, enhancing availability. |
| Privacy Concerns | Ensures data confidentiality by creating anonymized synthetic data. |
| High Costs of Data Collection | Reduces reliance on real data, minimizing costs associated with clinical trials. |
| Limited Diversity in Datasets | Simulates diverse clinical scenarios, including rare diseases and demographic groups. |
| Ethical Constraints | Eliminates the need for direct patient involvement, reducing ethical dilemmas. |
| Time-Consuming Data Access and Sharing | Facilitates faster data sharing by creating readily available synthetic datasets. |
| Bias in Data Representation | Generates balanced datasets to reduce biases inherent in real-world data. |
| Inflexibility in Modeling Rare Scenarios | Enables simulation of rare clinical conditions and complex disease progressions. |
| Regulatory Hurdles | Bypasses restrictions on patient data usage while adhering to privacy regulations. |
| Lack of Patient Representation in AI Training Data | Improves representation by synthesizing underrepresented patient profiles. |

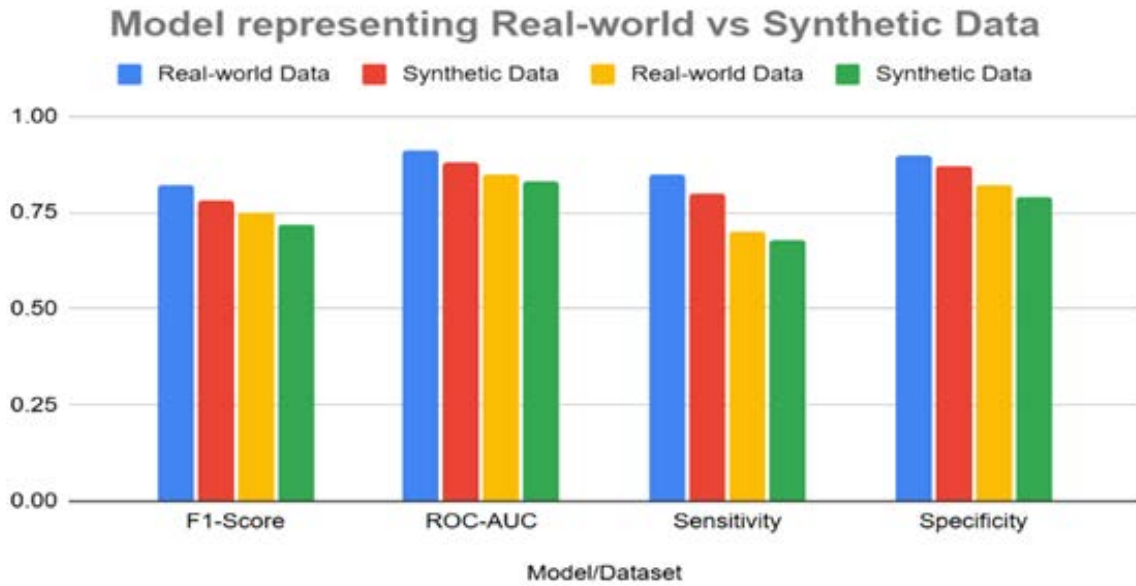**Table 2:** Comparison of Synthetic Data Generation Techniques.

| Method | Advantages | Limitations | Use Cases |
|---|---|---|---|
| Generative Adversarial Networks (GANs) | Produces highly realistic synthetic data. Can capture complex patterns in data. | Requires large datasets for training. Susceptible to mode collapse (lack of diversity). | Synthetic medical images (e.g., MRIs, X-rays). Disease progression simulation. |
| Variational Autoencoders (VAEs) | Generates interpretable and diverse data. Learns probabilistic data representations. | Output may lack realism compared to GANs. Training can be computationally intensive. | Synthetic patient records. Modeling disease co-occurrence and progression. |
| Gaussian Copula | Simple to implement. Preserves statistical relationships in data. | Struggles with high-dimensional or non-linear data. | Tabular data generation for clinical trials. |
| Tabular Variational Autoencoders (TVAEs) | Optimized for tabular data. Captures dependencies between variables effectively. | Requires parameter tuning. Potential loss of interpretability in some cases. | Synthetic EHRs for privacy-preserving research. |
| Diffusion Models | Generates high-fidelity data by iteratively refining noise. | Computationally expensive. Relatively new, with fewer established healthcare applications. | Medical image generation for training diagnostic AI models. |
| Multivariate Information Framework (MIIC-SDG) | Balances privacy and utility. Effective in creating realistic data distributions. | Limited adoption in large-scale datasets. | Longitudinal cohort data for predictive modeling. |
| ModOSIM | Closely resembles real-world administrative health data. | Applicability limited to specific healthcare systems. | Drug safety studies and analyst training. |

**Table 3:** Synthetic Data of breast cancer patients of 40-65 years of age belonging to Hispanic and non-hispanic ethnicity.

| Age | Gender (0=Female, 1=Male) | Race (1=White, 2=Black, 3=Asian, 4=Other) | Ethnicity (0=Non-Hispanic, 1=Hispanic) | Medical History (1=BRCA1, 2=BRCA2, 3=Family History, 4=No Significant History) | Tumor Stage (0=Early, 1=Late) | Treatment Plan (1=Surgery, 2=Chemotherapy, 3=Radiation, 4=Hormone Therapy, 5=Combination) |
|---|---|---|---|---|---|---|
| 55 | 0 | 1 | 0 | 3 | 1 | 5 |
| 42 | 0 | 2 | 0 | 1 | 0 | 1 |
| 60 | 0 | 3 | 1 | 4 | 1 | 2 |
| 48 | 0 | 1 | 0 | 2 | 0 | 3 |
| 52 | 0 | 4 | 0 | 3 | 1 | 5 |
| 58 | 0 | 2 | 1 | 4 | 0 | 1 |
| 45 | 0 | 1 | 0 | 2 | 1 | 5 |
| 62 | 0 | 3 | 0 | 3 | 1 | 2 |
| 50 | 0 | 4 | 1 | 4 | 0 | 3 |
| 53 | 0 | 1 | 0 | 1 | 1 | 5 |



**Figure 1:** Graphical Representation.

**Figure 2:** Graph representing the Performance Comparison of Diagnostic Models.

**Table4:** Specialists Evaluation of Synthetic Data Plausibility.

| Specialists | Credibility Rating (1-5) | Key Observations |
|---|---|---|
| Cardiologist | 4 | "The synthetic data appears realistic, but some edge cases and rare conditions may not be adequately represented." |
| Cardiologist | 3 | "While the overall data quality is good, there are inconsistencies in certain patient histories and lab results." |
| Cardiologist | 5 | "The synthetic data is highly credible and indistinguishable from real-world data. It could be used for various clinical research purposes." |
| Oncologist | 4 | "The data generation process needs refinement to improve the accuracy of certain demographic and clinical features." |
| Oncologist | 3 | "The synthetic data is useful for training AI models, but it may not be suitable for complex clinical decision-making scenarios." |
| Oncologist | 4 | "The synthetic data is generally credible, but there are some limitations in terms of capturing the complexity of real-world clinical variability." |
| Radiologist | 5 | "The synthetic data is of excellent quality and can be used to address a wide range of research questions." |
| Radiologist | 3 | "While the data is realistic, it lacks the depth and nuance of real-world clinical data, particularly in terms of rare diseases and comorbidities." |
| Radiologist | 4 | "The synthetic data is a valuable tool for training and testing AI models, but it should be used in conjunction with real-world data for validation." |
| Data Scientist | 2 | "The synthetic data has significant limitations in terms of data quality and clinical relevance. Further improvements are needed to make it a reliable source of information." |

**Table 5:** Ethical Risk vs. Benefit Analysis of Synthetic Data Use.

| Ethical Risk | Potential Benefit | Mitigation Strategies |
|---|---|---|
| Privacy Concerns | Improved Data Accessibility | Strong anonymization and de-identification techniques<br>Differential privacy methods<br>Regular privacy impact assessments |
| Bias and Fairness | Reduced Bias in AI Models | Careful selection and curation of training data<br>Bias detection and mitigation techniques<br>Regular model evaluation for fairness |
| Misuse and Misinterpretation | Enhanced Research and Innovation | Clear guidelines and regulations for synthetic data use<br>Education and training for researchers and practitioners<br>Transparent documentation of data generation and limitations |
| Data Quality and Realism | More Robust AI Models | Rigorous evaluation of synthetic data quality<br>Continuous improvement of data generation techniques<br>Regular validation against real-world data |
| Legal and Regulatory Challenges | Accelerated Drug Discovery and Clinical Trials | Proactive engagement with regulatory bodies<br>Development of clear legal frameworks for synthetic data |

## 6. Conclusion

The integration of synthetic data through generative AI represents a transformative breakthrough in clinical research and healthcare. By addressing the long-standing challenges of data scarcity, privacy concerns and accessibility, generative models such

as GANs and VAEs have opened new doors to a more inclusive, precise and equitable medical landscape. Synthetic data not only preserves the statistical integrity of real-world information but also protects patient confidentiality, fostering trust and ethical responsibility in an era where data is paramount.

Findings in this paper underscore the profound potential of synthetic data in revolutionizing diagnostic accuracy, enhancing clinical trial designs and driving the development of personalized medicine. The ability to simulate diverse patient scenarios empowers researchers to predict disease trajectories and optimize treatment strategies, ultimately bridging gaps in knowledge and resources that have historically constrained innovation. The promise of generative AI is evident in its capacity to accelerate drug discovery, simulate rare conditions and expand the horizons of medical insights while mitigating risks of data misuse and ethical breaches.

However, this progress comes with responsibility. The limitations of AI models, including inaccuracies and biases, highlight the necessity for continuous refinement, rigorous validation and strong regulatory frameworks. The insights of clinical experts, paired with advanced ethical analysis, remind us that the journey toward integrating synthetic data must be guided by a steadfast commitment to patient welfare, transparency and accountability.

This study envisions a future where every step forward in clinical research is a testament to innovation that is as compassionate as it is groundbreaking. Generative AI, when wielded responsibly, has the potential to redefine healthcare, ensuring that better diagnostics, effective treatments and equitable care are no longer privileges but rights for all. This convergence of technology and humanity offers hope for a world where every patient's story is understood, every life valued and every medical milestone celebrated as a collective achievement.

## 7. References

1. Topol EJ. High-performance medicine: The convergence of human and artificial intelligence. Nature Medicine, 2019;25:44-56.

2. Obermeyer Z, Lee TH. Lost in thought - The limits of the human mind and the future of medicine. New England Journal of Medicine, 2017;377:1209-1211.

3. Esteva A, Robicquet A, Ramsundar B, Kuleshov V, DePristo M, Chou K, Dean J. A guide to deep learning in healthcare. Nature Medicine, 2019;25:24-29.

4. Liu X, Faes L, Kale AU, Wagner SK, Fu DJ, Bruynseels A, Maetschke S. A comparison of deep learning performance against health-care professionals in detecting diseases from medical imaging: A systematic review and meta-analysis. The Lancet Digital Health, 2020;2:271-297.

5. Choi E, Schuetz A, Stewart WF, Sun J. Using recurrent neural networks for early detection of heart failure from clinical data. Journal of the American Medical Informatics Association, 2017;24:266-272.

6. Frid-Adar M, Elter M, Kahn CE. GANs for generating synthetic medical images: A survey. IEEE Transactions on Biomedical Engineering, 2018;65:2152-2164.

7. Yang Y, Wei C, Xie Y. Synthetic patient data for drug discovery and clinical trials. Trends in Pharmacological Sciences, 2019;40:360-368.

8. Obermeyer Z, Emanuel EJ. Predicting the future - big data, machine learning and clinical medicine. New England Journal of Medicine, 2016;375:1216-1219.

9. Rajkomar A, Dean J, Kohane I. Machine learning in medicine. New England Journal of Medicine, 2018;380:1347-1358.

10. Desai S, Ohno-Machado L, Gombar S. Predicting sepsis in the ICU using deep learning. Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2017;1195-1204.

11. Zhang Y, Zhao X, Li J. Automated systematic review using natural language processing, 2020.

12. Char DS, Shah NH, Magnus D. Implementing machine learning in health care - Addressing ethical challenges. New England Journal of Medicine, 2018;378:981-983.

13. Reddy CK, Kambhampati C, Hasan MA. A survey of data privacy and security issues in electronic health records. IEEE Access, 2019;7:83243-83268.

14. Sweeney L, Langer S. Data quality in electronic health records: A review. Journal of Biomedical Informatics, 2019;94:103197.

15. Goodfellow I, Pouget-Abadie J, Mirza M. Generative adversarial nets. Proceedings of the 27th International Conference on Neural Information Processing Systems, 2014;2672-2680.

16. Sim JJM, Rusli KDB, Seah B, Levett-Jones T, Liaw SY. Virtual simulation to enhance clinical reasoning in nursing: A systematic review and meta-analysis. Clinical Simulation in Nursing, 2022;69:26-39.

17. Nie D, Wang L, Yang X. Medical image synthesis with deep learning methods. Journal of Medical Imaging, 2020;7:1-22.

18. Zhang L, Zhang L, Wang J. Generative models for drug discovery. Nature Reviews Drug Discovery, 2020;19:489-510.

19. Wang Q, Chen J, Zhang X. Generative models for personalized medicine: A review. Frontiers in Genetics, 2021;12:649835.

20. Howe Iii EG, Elenberg F. Ethical Challenges Posed by Big Data. Innov Clin Neurosci, 2020;17:24-30.

21. Kingma DP, Welling M. Auto-Encoding Variational Bayes, 2014.

22. Rezende DJ, Mohamed S. Variational Inference with Normalizing Flows. Proceedings of the 32nd International Conference on Machine Learning, 2015;1530-1538.

23. Frid-Adar M, Diamant I, Klang E, et al. GAN-based synthetic medical image augmentation for increased CNN performance in liver lesion classification. Neurocomputing, 2018;321:321-331.

24. Chandrakant Mallick, Parimal Kumar Giri, Bijay Paikaray. The Privacy-Preserving High-Dimensional Synthetic Data Generation and Evaluation in the Healthcare Domain. Advances in data mining and database management book series, 2024.

25. Richard K Lomotey, Sandra Kumi, Madhurima Ray, Ralph Deters. Synthetic Data Digital Twins and Data Trusts Control for Privacy in Health Data Sharing, 2024.

26. Jennifer Anne Bartell, Sander Boisen Valentin anders Krogh, Henning Langberg, Martin Bøgsted. A primer on synthetic health data, 2024.

27. Nadir Sella, Florent Guinot, Nikita Lagrange, Laurent-Philippe Albou, Jonathan Desponds, Hervé Isambert. Preserving Information while Respecting Privacy: An Information Theoretic Framework for Synthetic Health Data Generation, 2024.

28. Olawale F Ayilara, Robert W, Platt, Matt Dahl, Janie Coulombe, Pablo Gonzalez Ginestet, Dan Château, Lisa M Lix. Generating synthetic data from administrative health records for drug safety and effectiveness studies. International Journal for Population Data Science, 2023.

29. Lisa Langnickel, John H Schneider, Ines Perrar, Tim Adams, Sobhan Moazemi, Fabian Praßer, Ute Nöthlings, Holger Fröhlich, Juliane Fluck. Synthetic data generation for a longitudinal cohort study - evaluation, method extension and reproduction of published data analysis results. Dental science reports, 2024.

30. Elnaz Karimian Sichani, Aaron Smith, Khaled El Emam, Lucy Mosquera. Creating High-Quality Synthetic Health Data: Framework for Model Development and Validation. JMIR formative research, 2023.

31. Zhaozhi Qian, Bogdan-Constantin Cebere S Janes, Neal Navani, Mihaela van der Schaar. Synthetic data for privacy-preserving clinical risk prediction, 2023.

32. Mohd Rafatullah. Synthetic data: the future of open-access health-care datasets? The Lancet, 2023.

33. Synthetic Data and amp; the Future of Women's Health: A Synergistic Relationship, 2023.