


# Chapter 6

## Generative AI Approaches for Synthetic Health Data Generation: Methods, Applications, and Challenges

V. Vanitha

 <https://orcid.org/0000-0003-1010-3241>

*Sri Ramachandra Faculty of Engineering and Technology, India*

M. Srivani

*Sri Ramachandra Faculty of Engineering and Technology, India*

### ABSTRACT

*Generating realistic synthetic health data using Artificial Intelligence models offers a powerful solution to privacy restrictions and data scarcity in healthcare. Modern generative frameworks such as Generative Adversarial Networks, Variational Autoencoders, diffusion models, and transformer-based language models enable the synthesis of diverse clinical data, including structured EHR tables, time-series vitals, unstructured clinical notes, and medical images. These models are evaluated on fidelity, utility, and privacy, with applications ranging from data augmentation to privacy-preserving data sharing. GANs excel at learning joint distributions, VAEs capture smooth latent structures, diffusion models yield high-fidelity samples, and transformers produce coherent medical text. Challenges remain in ensuring clinical validity, diversity, and robust privacy protection while avoiding biases and mode collapse. Emerging research explores multimodal generation, promising richer, more integrated synthetic patient records that can transform healthcare research and AI model development.*

DOI: 10.4018/979-8-3373-5641-9.ch006

# 1. INTRODUCTION

Synthetic health data have surged in importance due to privacy regulations and the high cost of collecting large medical datasets. Modern healthcare relies on large-scale electronic health records (EHRs), imaging archives, and genomic databases, but sharing these data is often restricted by laws (e.g., HIPAA, GDPR) or institutional policies. At the same time, many medical applications—such as developing AI models for disease prediction—require abundant, diverse, and richly annotated data.

By mimicking the statistical properties of real data, synthetic datasets enable broader data sharing, data augmentation, and fairness in analysis. Ibrahim et al. (2025) note that synthetic data facilitate sharing while protecting privacy, augment existing datasets, and promote fairness in AI applications. Similarly, Goncalves et al. (2020) emphasize that the scarcity of patient data has slowed progress in medicine, making high-quality synthetic data critical when real data are difficult, expensive, or limited to acquire.

Synthetic datasets can help address the lack of representative samples (e.g., for rare diseases) and accelerate the development of diagnostic models and statistical studies. By mirroring the statistical structures of real data while severing links to individuals, they allow researchers to access realistic datasets without exposing actual patient information. In practice, carefully generated synthetic data support broad sharing, foster algorithm development, and serve as benchmarks—all with minimal privacy risk.

## 1.1 The Importance of Synthetic Health Data

The key benefits of synthetic health data include:

1. **Data augmentation** – Expanding scarce datasets, especially for rare conditions or under-represented populations. For instance, synthetic images of rare tumors or EHR records of uncommon comorbidities can improve model training and reduce overfitting. Techniques such as Variational Autoencoders and GANs have successfully generated medical images (e.g., pneumonia in chest X-rays) and clinical signals, improving performance on minority classes.
2. **Privacy preservation** – High-quality synthetic EHRs can be shared between institutions without leaking Protected Health Information (PHI). The records mimic real data but are decoupled from real individuals, minimizing re-identification risk, (Nikolentzos et al., 2023).
3. **Fairness promotion** – By oversampling or simulating data for vulnerable or under-represented groups, synthetic datasets can help mitigate bias and improve the fairness of AI models.

30 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: [www.igi-global.com/chapter/generative-ai-approaches-for-synthetic-health-data-generation/397163](http://www.igi-global.com/chapter/generative-ai-approaches-for-synthetic-health-data-generation/397163)

## Related Content

---

### Productivity Profiles of Islamic Banks Using Data Envelopment Analysis-Based Malmquist Productivity Indices (MPIs): Survey, Classification, and Critical Analysis

Karim Iddouch, Khalid El Badraoui and Jamal Ouenniche (2024). *Data Envelopment Analysis (DEA) Methods for Maximizing Efficiency* (pp. 40-81).

[www.irma-international.org/chapter/productivity-profiles-of-islamic-banks-using-data-envelopment-analysis-based-malmquist-productivity-indices-mpis/336940](http://www.irma-international.org/chapter/productivity-profiles-of-islamic-banks-using-data-envelopment-analysis-based-malmquist-productivity-indices-mpis/336940)

### Big Data Analysis: Basic Review on Techniques

Arpit Kumar Sharma, Arvind Dhaka, Amita Nandal, Kumar Swastik and Sunita Kumari (2021). *Advancing the Power of Learning Analytics and Big Data in Education* (pp. 208-233).

[www.irma-international.org/chapter/big-data-analysis/272956](http://www.irma-international.org/chapter/big-data-analysis/272956)

### Commercial Banks' Digital Paradigm and Customers Responses in the UAE

Muhammad Jumaa (2020). *International Journal of Data Analytics* (pp. 68-79).

[www.irma-international.org/article/commercial-banks-digital-paradigm-and-customers-responses-in-the-uae/244170](http://www.irma-international.org/article/commercial-banks-digital-paradigm-and-customers-responses-in-the-uae/244170)

### Big Data Analytics With Machine Learning and Deep Learning Methods for Detection of Anomalies in Network Traffic

Valliammal Narayanand Shanmugapriya D. (2022). *Research Anthology on Big Data Analytics, Architectures, and Applications* (pp. 678-707).

[www.irma-international.org/chapter/big-data-analytics-with-machine-learning-and-deep-learning-methods-for-detection-of-anomalies-in-network-traffic/291007](http://www.irma-international.org/chapter/big-data-analytics-with-machine-learning-and-deep-learning-methods-for-detection-of-anomalies-in-network-traffic/291007)

### ICTs and Domestic Violence (DV): Exploring Intimate Partner Violence (IPV)

Bolanle A. Olaniran (2021). *International Journal of Big Data and Analytics in Healthcare* (pp. 31-44).

[www.irma-international.org/article/icts-and-domestic-violence-dv/277646](http://www.irma-international.org/article/icts-and-domestic-violence-dv/277646)