

# Navigating Opportunities and Challenges in Synthetic Data Generation for Biomedicine: Insights from the SYNTHIA Project



Patricia Alonso de Apellániz <sup>1</sup>, Aleksandar Babic <sup>2</sup>, Vibeke Binz Vallevik <sup>2,3</sup>, Oscar Brück <sup>4</sup>, Gastone Castellani <sup>5</sup>, Chiara Chianese <sup>6</sup>, Davide Cirillo <sup>7</sup>, Vanda Czipczer <sup>8</sup>, Saverio D’Amico <sup>9,10</sup>, Mattia Delleani <sup>9</sup>, Matteo Della Porta <sup>10</sup>, Nicolas Derus <sup>5</sup>, Chong Duan <sup>11</sup>, Styliani-Christina Fragkouli <sup>12</sup>, Holger Fröhlich <sup>13</sup>, Adrian Galiana-Bordera <sup>14</sup>, Holger Hennig <sup>15</sup>, Katja Herzog <sup>15</sup>, Ida Holte Thorius <sup>6</sup>, Bayrem Kaabachi <sup>16</sup>, Jan Klein <sup>17</sup>, Bogdan Kulynych <sup>16</sup>, Levente Lippenszky <sup>8</sup>, Luis Marti-Bonmati <sup>14</sup>, Serena Elizabeth Marshall <sup>2</sup>, Nagat Masued <sup>7</sup>, Daniel Mensing <sup>17</sup>, Matthias Mullenborn <sup>6</sup>, Juan Parras Moral <sup>1</sup>, Fotis Psomopoulos <sup>12</sup>, Núria Queralt-Rosinach <sup>18</sup>, Jean Louis Raisaro <sup>16</sup>, Diego Valderrama <sup>13</sup>, Hongxu Yang <sup>8</sup>, Ahmed Youssef Ali <sup>19</sup>, Santiago Zazo <sup>1</sup>

<sup>1</sup> Universidad Politécnica de Madrid (UPM), <sup>2</sup> DNV, <sup>3</sup> University of Oslo (UiO), <sup>4</sup> Helsinki University Hospital (HUS), <sup>5</sup> Università di Bologna (UNIBO), <sup>6</sup> Novo Nordisk, <sup>7</sup> Barcelona Supercomputing Center (BSC), <sup>8</sup> GE Healthcare, <sup>9</sup> TRAIN, <sup>10</sup> Humanitas Research Hospital, <sup>11</sup> Pfizer, <sup>12</sup> Centre for Research and Technology-Hellas (CERTH), <sup>13</sup> Fraunhofer SCAI, <sup>14</sup> Instituto de Investigación Sanitaria La Fe (IIS La Fe), <sup>15</sup> Fraunhofer ITMP, <sup>16</sup> Centre hospitalier universitaire vaudois (CHUV), <sup>17</sup> Fraunhofer MEVIS, <sup>18</sup> Leiden University Medical Center (LUMC), <sup>19</sup> Johnson & Johnson



Find Out More!

## The IHI SYNTHIA project

Funded by the **Innovative Health Initiative (IHI)**, SYNTHIA is a public-private partnership advancing **synthetic data generation (SDG)** for biomedicine. It develops tools to tackle data scarcity and privacy, generating tabular, imaging, sequencing data and more across six disease areas: **lung cancer, breast cancer, multiple myeloma, DLBCL, Alzheimer’s, and Type 2 diabetes** (Fig. 1). Outputs will be delivered via a **federated platform** to support research and innovation.

## A Comprehensive Review of SDG Methods in Biomedicine

As part of the SYNTHIA project, a **scoping review** was conducted to assess the state of the art, identify key challenges and opportunities, and define a shared understanding of SDG to guide future development. The review followed a structured methodology, focusing on high-impact publications and preprints from the past 5–10 years across databases including **PubMed, Scopus, Web of Science, IEEE Xplore, and arXiv**.

## Overview of common SDG approaches

| Approach         | Strengths                  | Limitations                            |
|------------------|----------------------------|--|
| Statistical      | Simple, interpretable      | Privacy, scalability, limited dynamics |
| Machine learning | Flexible, powerful         | Complex, compute-heavy, training risks |
| Simulations      | Domain-grounded, efficient | Less adaptable, lower variability      |

### Statistical methods

- **KDE:** Non-parametric, replicates distributions; limited privacy, high compute.
- **Gaussian Copulas:** Capture multivariate dependencies; struggle with scale/privacy.
- **Mixture Models (e.g., GMMs):** Model heterogeneity; newer versions handle mixed data.
- **Bayesian Networks:** Encode dependencies; some support private variants.
- **Oversampling (SMOTE/ADASYN):** Balance classes; simple but privacy-limited.

### Machine learning methods

- **NODEs / ANODEs:** Model latent continuous-time dynamics; great for irregular time-series.
- **Neural Laplace:** Laplace-enhanced NODEs; better for sparse time points.
- **GNNs:** Learn on graph data; with variants for relational health data.
- **VAEs:** Probabilistic encoding; some handle missing/mixed data and temporal modeling.
- **GANs:** Realistic outputs via adversarial learning; unstable, mode collapse risk.
- **DDPMs:** Probabilistic noise-based generation; high-quality, but slow.
- **LLMs (e.g., GPTs):** Emerging for tabular/text; expressive but with privacy/ethics concerns.

### Simulative methods

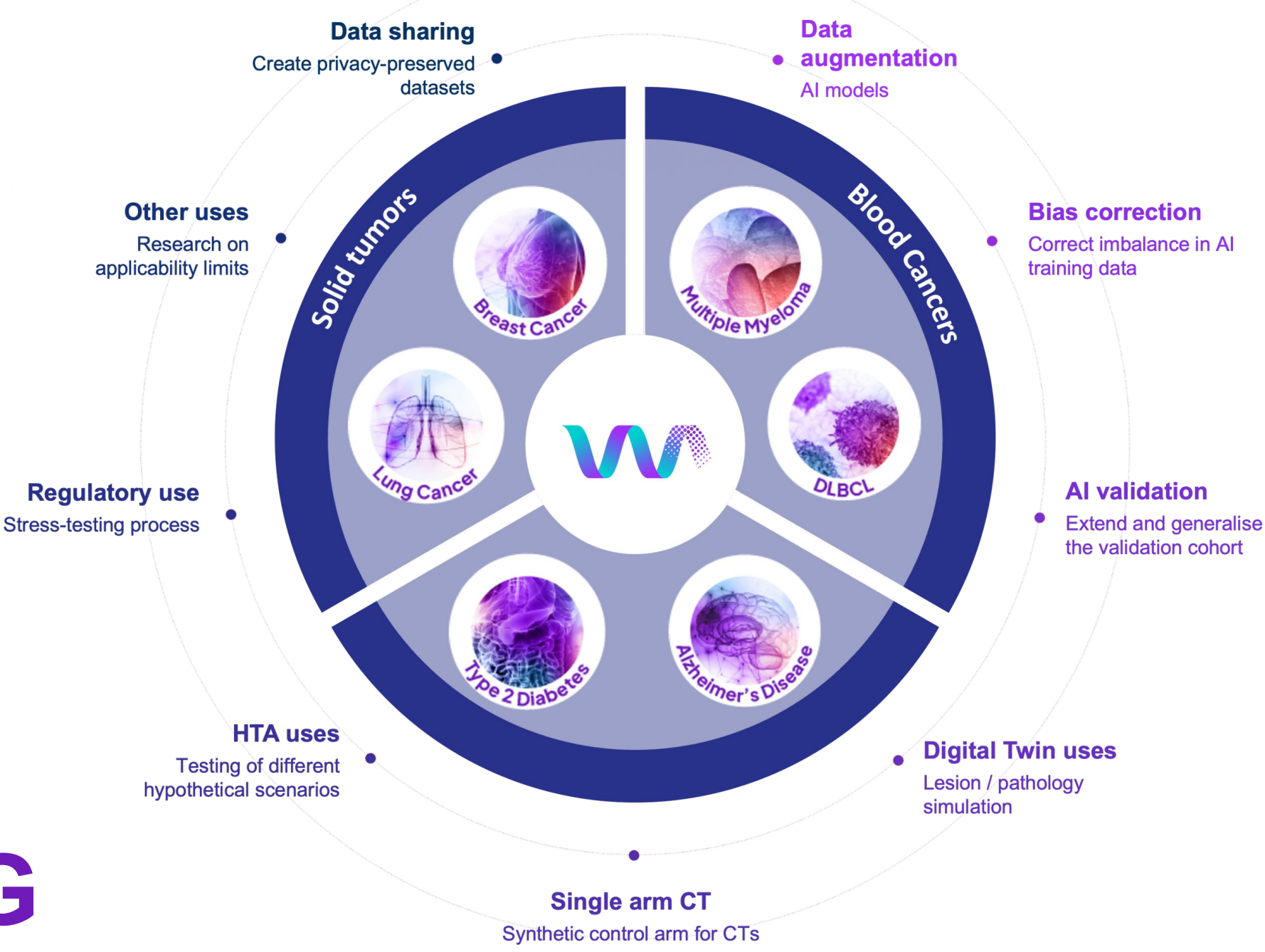
- **Imaging-Based:** Embed clinical features in clean images; efficient and realistic.
- **Physics/Chemistry Models:** Simulate interactions (e.g., contrast imaging).
- **Rule-Based Simulators (e.g., Synthea):** Generate synthetic EHRs from clinical rules.

### Multimodal data

Multimodal synthetic data generation in healthcare leverages advanced AI techniques to integrate diverse data types, such as tabular, imaging, time-series, and omics, within unified frameworks. Methods include, transformer-based architectures, hybrid models, multimodal GANs, and diffusion models, which learn complex cross-modal dependencies.

## CONCLUSIONS

Ensuring data quality, clinical relevance, and minimizing bias requires robust **assessment, benchmarking** and **FAIR principles**. The evolving regulatory landscape (e.g., GDPR, HIPAA) calls for clear legal and ethical guidance. Future efforts must focus on **standardized evaluation metrics, stronger privacy protections, and expert-in-the-loop validation** to ensure real-world clinical utility.



### Synthetic Data Generation framework for integrated validation of use cases and AI healthcare applications.

This project is supported by the Innovative Health Initiative Joint Undertaking (IHI JU) under grant agreement No 101172872. The JU receives support from the European Union's Horizon Europe research and innovation programme, COCIR, EFPIA, Europa Bio, MedTech Europe, Vaccines Europe and DNV. The UK consortium partner, The National Institute for Health and Care Excellence (NICE) is supported by UKRI Grant 10132181.



Funded by the European Union, the private members, and those contributing partners of the IHI JU. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the aforementioned parties. Neither of the aforementioned parties can be held responsible for them.

**Figure 1. Targeted Synthetic Data Applications.** SYNTHIA delivers purpose-built synthetic datasets and tools to meet specific clinical and research needs in real-world healthcare settings.

## SDG methods by data modality

### Textual data

Synthetic medical text is generated using LLMs like GPT-4 and transformer-based frameworks such as **MedSyn** [1], which integrates Medical Knowledge Graphs, improving NLP tasks like ICD coding, NER, and de-identification. Tools like DeID-GPT automate removal of private information, enabling realistic, privacy-safe clinical text generation for research and clinical use.

### Tabular data

Synthetic tabular data in healthcare is generated using classical statistical methods like GMMs, Copulas, and KDE, as well as machine learning techniques including decision trees (e.g., **Randomized Decision Trees** [2]), GANs (e.g., **MedGAN** [3], **CTGAN** [4]), VAEs (e.g., **TVAE** [4], **TabVAE** [5]), diffusion models (e.g., **TabDDPM** [6], **TabSyn** [7]), and LLMs (e.g., **GReaT** [8]). These methods address challenges like mixed data types, privacy, and complex dependencies, supporting applications from clinical trials to precision medicine.

### Imaging data

Medical imaging data synthesis leverages tools like GANs, diffusion models, and hybrid approaches (e.g., **HA-GAN** [9]) to generate anatomically accurate and scalable synthetic images while addressing data scarcity and privacy concerns. Additionally, vision-language models such as **MedViLL** [10] and **Flamingo-CXR** [11] enable automated generation of clinically relevant radiology reports, improving workflow efficiency and diagnostic support.

### Times series data

Synthetic signaling and time-series data generation utilizes advanced tools like **VAMBN** [12] and **VAMBN-MT** [13] for capturing complex temporal and multimodal dependencies, as well as **MultiNODEs** [14] for modeling continuous trajectories in latent space. For biomedical signals such as EEG and ECG, GAN variants (e.g., **CGANs** [15], **RGANs** [16]), **DDPMs** [17], **SynSigGAN** [18], and **DoppelGANger** [19] are employed to create realistic, high-fidelity synthetic data reflecting temporal dynamics and signal complexity.

[1] Kurnichev G, Bilinov P, Kuzkina Y, et al. *Lect Notes Comput Sci*. 2024;215–30.  
[2] Vaidya J, Shafiq B, Asani M, et al. *AMIA Annu Symp Proc*. 2017;1695–704.  
[3] Choi E, Biswal S, Malin B, et al. *arXiv*. 2017.  
[4] Xu L, Skoulikidou M, Quesada-Trifante A, et al. *arXiv*. 2019.  
[5] Tazawa S, Knobloch M, Quesada E, et al. *Proc ICPRAM*. 2024.  
[6] Kotelnikov A, Baranchuk D, Rubachev I, et al. *arXiv*. 2022.  
[7] Zhang H, Zhang J, Srinivasan B, et al. *arXiv*. 2023.  
[8] Borisov V, Seifler K, Leemann T, et al. *arXiv*. 2022.  
[9] Sun L, Chen J, Xu Y, et al. *IEEE J Biomed Health Inform*. 2022;26:3966–75.  
[10] Ye-Bin M, Hyeon-Woo N, Choi W, et al. *arXiv*. 2023.  
[11] Tanno R, Barrett DT, Seifried A, et al. *Nat Med*. 2024.  
[12] Goodes-Dresbach L, Sood M, Sahay A, et al. *Front Big Data*. 2020;3:16.  
[13] Kühnel L, Schneider J, Perrar I, et al. *Sci Rep*. 2024;14:14412.  
[14] Wendland P, Birkenbihl C, Gomez-Freixa M, et al. *NPJ Digit Med*. 2022;5:122.  
[15] Wickramaratne SD, Parekh A. *Conf IEEE EMBC*. 2023;1–4.  
[16] Esteban C, Hyland SL, Ratsch G. *arXiv*. 2017.  
[17] Ho J, Jain A, Abbeel P. *arXiv*. 2020.  
[18] Hazza D, Byun Y-G. *Biology (Basel)*. 2020;9.  
[19] Lin Z, Jain A, Wang C, et al. *Proc ACM IMC*. 2020.  
[20] Nguyen E, Poli M, Durrant MG, et al. *bioRxiv*. 2024.  
[21] Zvyagin M, Brace A, Hippe K, et al. *bioRxiv*. 2022.  
[22] Chen B, Cheng X, Li P, et al. *bioRxiv*. 2023.  
[23] Song D, Wang Q, Yan G, et al. *Nat Biotechnol*. 2024;42:247–52.