

AtmosGen: Condition-Aware Synthetic Atmospheric Data and Image Generation for Aviation Applications

Pradeep Shirvi¹, Srushti Patil¹, and Aakashi Jangam¹

¹Students, Department of Artificial Intelligence and Machine Learning, Universal College of Engineering, Vasai, India

pradeepshirvi72@gmail.com, Srushti08patil@gmail.com, Aakashi0904@gmail.com,

Abstract

Abstract – Accurate atmospheric data is essential for weather forecasting, aviation safety, and climate research. However, real-world data collection methods such as radiosonde launches are limited by high operational costs, sparse temporal availability, and restricted geographical coverage. To address these challenges, this paper proposes AtmosGen (Atmospheric Synthetic Data and Image Generator), a condition-aware synthetic atmospheric data and image generation framework that combines numerical data synthesis with atmospheric image generation. The system utilizes historical radiosonde data to generate realistic synthetic atmospheric parameters, including temperature, pressure, humidity, wind speed, and altitude, using machine learning-based generative models. In addition, a conditional image generation model is employed to generate atmospheric and weather-condition images corresponding to different environmental states such as clear sky, cloudy, foggy, and stormy conditions. To ensure the reliability of the generated data, a compatibility evaluation model is introduced, which verifies the consistency between input atmospheric conditions and the generated images using statistical similarity metrics and regression-based validation. Furthermore, a comparative analysis between original radiosonde datasets and model-generated datasets is performed using distribution analysis, correlation metrics, and downstream task performance evaluation. The proposed approach reduces dependency on continuous real-time data acquisition while providing scalable, diverse, and scientifically consistent datasets. This framework is particularly useful for aviation simulations, machine learning model training, and atmospheric research where large labeled datasets are required.

Keywords: Synthetic Atmospheric Data, Image Generation, Radiosonde, Conditional GAN, Aviation Weather Simulation

I. Introduction

Accurate atmospheric data is essential for aviation safety, weather forecasting, and climate research. Aircraft operations rely on vertical atmospheric profiles such as temperature, pressure, humidity, and wind speed, typically obtained from radiosonde observations. Although radiosondes provide reliable upper-air measurements, they are limited

by high operational costs, fixed launch schedules, sparse geographical coverage, and discontinuous data availability. These limitations create data gaps that restrict large-scale simulations and machine learning applications requiring continuous atmospheric datasets.

Recent advancements in machine learning and generative modeling have enabled the development of synthetic data generation systems capable of replicating real-world statistical patterns. Synthetic atmospheric datasets can reduce dependency on costly balloon launches while supporting scalable research and aviation analytics. However, generating scientifically consistent synthetic atmospheric data that preserves physical realism and statistical similarity to real observations remains a significant challenge.

This paper proposes AtmosGen, a Synthetic Atmospheric and Aviation Dataset Generator using historical radiosonde data as the reference source. The system applies machine learning-based generative models to produce realistic atmospheric profiles, including temperature, pressure, humidity, and wind speed across altitude levels. A compatibility and comparison model evaluates statistical similarity between original and generated datasets using correlation metrics and error analysis. The framework also supports condition-aware mapping between generated atmospheric parameters and corresponding weather representations.

Project Overview: Design and Functionality

The proposed system is designed as a structured data-driven atmospheric intelligence system that transforms raw radiosonde observations into validated synthetic datasets through a layered data processing architecture. The core components are described below.

Core Focus – Synthetic Atmospheric Profile Generation

The primary objective is to generate realistic atmospheric vertical profiles using historical radiosonde data, simulating key parameters such as:

- Temperature
- Pressure
- Humidity
- Wind Speed

Supporting Parameters – Aviation & Environmental Factors

To enhance realism, the system integrates contextual aviation-related parameters, including:

- Turbulence indicators
- Wind shear patterns
- Seasonal variation
- Sensor noise simulation

Predictive & Generative Modeling

The system employs the following generative models:

- Generative Adversarial Networks (GAN)
- Variational Autoencoders (VAE)
- Statistical regression models

A comparison model validates outputs against real data using Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and Pearson Correlation Coefficient.

Visualization & Decision Support

Outputs include synthetic dataset previews, statistical comparison results, correlation graphs, and atmospheric trend visualizations in CSV/JSON formats to support downstream research and simulation tasks.

Existing Systems and Their Limitations

Traditional atmospheric data collection systems primarily rely on radiosonde-based observations for obtaining upper-air measurements such as temperature, pressure, humidity, and wind speed at various altitude levels. Radiosondes, which are balloon-borne instruments, have long served as a reliable source of vertical atmospheric profiling for aviation safety, weather forecasting, and climate research. These observations provide high-quality in-situ measurements and are widely used in meteorological modeling and aviation planning. In addition to radiosondes, satellite observations and numerical reanalysis models are also used to supplement atmospheric data collection. However, radiosonde launches are typically conducted only once or twice daily from fixed geographical stations, leading to sparse temporal coverage and limited real-time adaptability. The operational process involves considerable costs related to balloon equipment, sensors, and manpower, making continuous large-scale deployment economically challenging. Furthermore, many remote regions and oceanic areas lack radiosonde stations, resulting in geographical data gaps. While satellite systems provide broader coverage, they often lack the high-resolution vertical profiling accuracy offered by radiosondes. Numerical reanalysis models depend heavily on available observational inputs, and inaccuracies in input data can propagate through the modeling system. Another limitation is the dependence on real-time physical data collection, which restricts scalability for applications that require large, continuous, and diverse datasets. The absence of synthetic augmentation mechanisms limits the ability to simulate rare atmospheric events such as turbulence spikes, extreme wind shear, or sudden pressure changes. Moreover, traditional systems do not provide built-in statistical compatibility validation when integrating generated or interpolated data, which may lead to inconsistencies in advanced modeling environments.

Problem Statement and Objective

Accurate, continuous, and high-resolution atmospheric data is essential for aviation safety, numerical weather prediction, climate monitoring, and data-driven machine learning applications. While radiosonde data provides reliable measurements, its availability

is constrained by high operational costs, limited launch frequency, sparse geographical distribution, and discontinuous temporal coverage, creating significant spatial and temporal data gaps especially in remote and oceanic regions. Existing alternatives such as satellite data and numerical reanalysis models provide broader coverage but often lack the vertical resolution accuracy offered by radiosondes. Additionally, current systems lack an integrated mechanism to generate statistically validated synthetic datasets capable of preserving physical realism while supporting scalable data augmentation. The primary objective of this research is to develop AtmosGen, which leverages historical radiosonde observations to produce realistic and scalable synthetic atmospheric profiles while preserving statistical distributions and physical interdependencies. The system incorporates aviation-related contextual variables such as turbulence effects, wind shear patterns, and seasonal variability to enhance simulation realism. A key objective is the development of a compatibility and comparison model using MAE, RMSE, and correlation analysis to validate generated data against real observations. Ultimately, the proposed system aims to reduce dependency on costly real-time data collection and provide a reliable atmospheric data simulation framework for aviation and research applications.

II. Review of Literature

Recent advancements in atmospheric and environmental monitoring have demonstrated the importance of multi-parameter data integration for accurate modeling and analysis. Ma et al. (2025) introduced an advanced atmospheric anomaly detection approach using integrated temperature and environmental parameters, improving the accuracy and reliability of atmospheric condition analysis. Their study demonstrated that combining multiple atmospheric indicators enhances the precision of environmental modeling and supports more reliable atmospheric data interpretation, highlighting the importance of multi-variable atmospheric modeling in synthetic dataset generation systems.

Gidey and Mhangara (2025) analyzed long-term atmospheric and environmental temperature relationships using satellite-based datasets and demonstrated strong correlations between atmospheric conditions and surface temperature variations. Their study emphasized the importance of accurate atmospheric data modeling for predictive analysis and environmental monitoring, supporting the need for scalable data generation frameworks that preserve statistical relationships between atmospheric parameters.

Garai et al. (2022) investigated the relationship between rainfall, temperature, and environmental conditions using regression-based modeling and satellite datasets. Their results confirmed strong correlations between atmospheric variables, demonstrating that modeling systems must preserve statistical relationships between parameters such as temperature, pressure, and humidity to ensure realistic atmospheric simulation.

Londhe et al. (2023) studied atmospheric and environmental variability across different climatic zones and observed strong seasonal and regional variations in atmospheric parameters, demonstrating that modeling systems must account for temporal variability to improve prediction accuracy and simulation reliability.

Awasthi et al. (2023) examined climate sensitivity and atmospheric variability and found significant relationships between atmospheric conditions and environmental changes, highlighting the importance of reliable atmospheric modeling frameworks for predictive

analysis and climate monitoring.

Dash et al. (2024) conducted a long-term study on atmospheric variability and precipitation patterns across multiple regions, emphasizing the importance of predictive modeling and atmospheric simulation systems for improving climate monitoring and forecasting capabilities.

Rani and Kumar (2023) analyzed the relationship between atmospheric parameters and environmental pollutants and demonstrated the influence of temperature, wind speed, and humidity on environmental stability, underscoring the need for accurate atmospheric data modeling for environmental monitoring.

Ravuri et al. (2021) proposed a deep generative model for skillful precipitation nowcasting using radar data, demonstrating that generative adversarial networks can produce physically consistent meteorological outputs that outperform traditional deterministic forecasting approaches. Their work directly supports the use of GAN-based architectures for realistic atmospheric data synthesis.

Stengel et al. (2020) investigated the use of GANs for super-resolution of climate model outputs, showing that generative models can preserve physical constraints and statistical distributions of atmospheric variables even when trained on sparse observational data. This work provides methodological grounding for physics-aware generative modeling in atmospheric contexts.

Price et al. (2023) presented a hybrid neural network and general circulation model framework for atmospheric simulation that incorporates physical realism constraints directly into the learning objective. Their findings underscore the importance of coupling data-driven generative models with physical consistency checks, a principle adopted in the compatibility validation component of AtmosGen.

Project Scope

The scope of AtmosGen is the creation of a unified atmospheric intelligence system that integrates radiosonde observations, aviation-related environmental parameters, and machine learning-based generative modeling into a single reliable framework. The system focuses on generating realistic synthetic atmospheric datasets while ensuring statistical consistency, physical realism, and compatibility with real-world atmospheric conditions.

Historical atmospheric data from radiosonde archives and meteorological repositories are used as primary inputs. The system performs preprocessing steps including data cleaning, normalization, missing value interpolation, and feature standardization to ensure that all parameters are consistent and suitable for generative modeling. A compatibility validation framework evaluates the similarity between real and synthetic datasets using MAE, RMSE, and correlation analysis, incorporating anomaly detection and statistical consistency checks to ensure scientifically meaningful outputs. Generative models including GANs, VAEs, and regression-based models are trained on historical data to learn temporal and spatial atmospheric patterns, enabling simulation of diverse aviation and environmental scenarios including rare or extreme atmospheric conditions.

III. Methodology

Building the Synthetic Atmospheric and Aviation Dataset Generator System

The development of AtmosGen followed a structured, data-engineering-centric methodology to transform limited and fragmented radiosonde atmospheric observations into a unified, scalable, and predictive synthetic atmospheric data platform integrating generative modeling, compatibility validation, and visualization.

1. Problem Identification

The core issue addressed is the limited availability, sparse coverage, and high operational cost of real-world radiosonde atmospheric data. Radiosondes are often discontinuous, geographically restricted, and insufficient for large-scale simulation and machine learning training. The proposed system acts as an atmospheric data generation and validation platform, addressing the key question: “How can realistic atmospheric datasets be generated when real observations are limited, and how can their compatibility with real-world atmospheric behavior be ensured?”

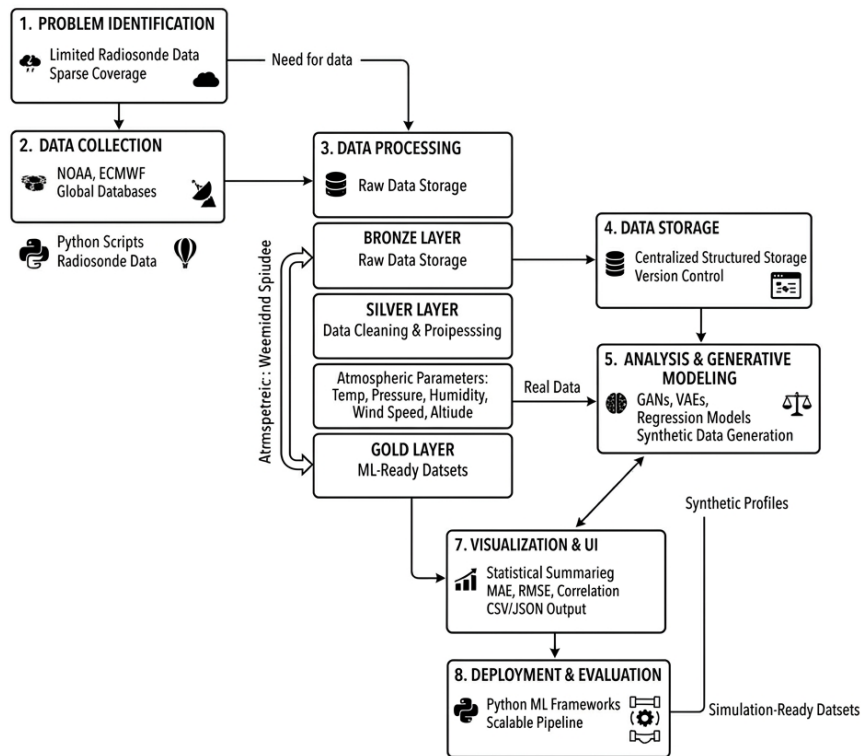


Figure 1: Block Diagram of AtmosGen

2. Data Collection

The system collects historical radiosonde atmospheric data from trusted meteorological repositories such as NOAA, ECMWF, and other global atmospheric databases. Python-based data extraction scripts retrieve and structure observations into tabular formats, stored in the initial data layer without modification to preserve integrity and traceability.

The training dataset comprises **50,000 atmospheric records** across **500 radiosonde-style vertical profiles**, each containing **100 altitude-level measurements** spanning 0 m to 30,000 m. Surface conditions were randomized to simulate diverse geographic and seasonal scenarios: surface temperature from -10°C to $+40^{\circ}\text{C}$, pressure from 990 hPa to 1040 hPa, and humidity from 30% to 95%. Each record contains seven features: altitude (m), temperature ($^{\circ}\text{C}$), pressure (hPa), relative humidity (%), wind speed (m/s), wind direction (degrees), and dewpoint ($^{\circ}\text{C}$). The dataset was generated using the International Standard Atmosphere (ISA) physics model with realistic stochastic variation, incorporating a lapse rate of -6.5°C per 1000 m, the barometric pressure formula, exponential humidity decay with altitude, and jet stream wind enhancement in the 8,000–12,000 m band. Preprocessing used StandardScaler normalization with fixed random seeds (NumPy seed 42, TensorFlow seed 42) for reproducibility. The full dataset was used for model training with no held-out split, as the evaluation framework relies on physics-model reference outputs for comparison.

3. Data Processing (The Transform Phase)

A structured layered data architecture ensures data quality and consistency. The **Bronze Layer** stores raw data in its original form. The **Silver Layer** performs data cleaning including removal of duplicate records, missing value interpolation, parameter normalization, and anomaly detection. The **Gold Layer** contains ML-ready processed datasets. Python libraries Pandas and NumPy are used throughout, with structured storage supporting efficient access and version control.

4. Data Storage and Warehousing

All processed datasets are stored in a centralized structured storage system supporting efficient querying, version control, and incremental data loading. This centralized storage serves as the primary source for model training, validation, and synthetic dataset generation.

5. Analysis and Generative Modeling

The system applies both descriptive and generative machine learning techniques to analyze atmospheric patterns and produce synthetic datasets. Descriptive analysis studies statistical distributions, correlations, and altitude-based variations. Two generative architectures were implemented, trained on the full 50,000-record dataset using a Google Colab T4 GPU and TensorFlow 2.x, and evaluated against the physics-based ISA baseline.

The **GAN (AtmosphericGAN)** generator comprises three fully connected hidden layers (256, 512, and 1024 neurons) with LeakyReLU activations ($\alpha = 0.2$) and Batch Normalization (momentum 0.8), accepting a 100-dimensional Gaussian noise vector as

input and producing a 7-dimensional atmospheric feature vector with Tanh activation. The discriminator uses layers of 512, 256, and 128 neurons with LeakyReLU activations and Dropout (rate 0.3), concluding with a Sigmoid output. Both networks were optimized using Adam (learning rate 0.0002, $\beta_1 = 0.5$) with binary cross-entropy loss for **2000 epochs** at batch size 64. The VAE (**TimeSeriesVAE**) encoder uses stacked LSTM layers (128 and 64 units) to project sequences of length 50 into a 20-dimensional latent space; the decoder reconstructs the 7-feature sequence via LSTM layers (64 and 128 units) with a RepeatVector bridge. Training used the Adam optimizer with combined MSE reconstruction and KL-divergence loss for **100 epochs** at batch size 32.

6. Visualization and User Interface

The system provides statistical summaries, correlation plots, and altitude-based atmospheric trend graphs through a structured visualization interface. Dataset outputs in CSV and JSON formats support downstream machine learning training, simulation, and analysis tasks.

7. Deployment and Evaluation

The system is deployed using Python-based machine learning frameworks and structured data pipelines. Model evaluation compares generated synthetic datasets against physics-model reference data using MAE, RMSE, and Pearson correlation. System evaluation covers data pipeline validation, model performance assessment, and dataset consistency verification.

IV. Result

The proposed AtmosGen system provides an interactive platform for generating, visualizing, and validating synthetic atmospheric and aviation datasets using configurable parameters and machine learning models.

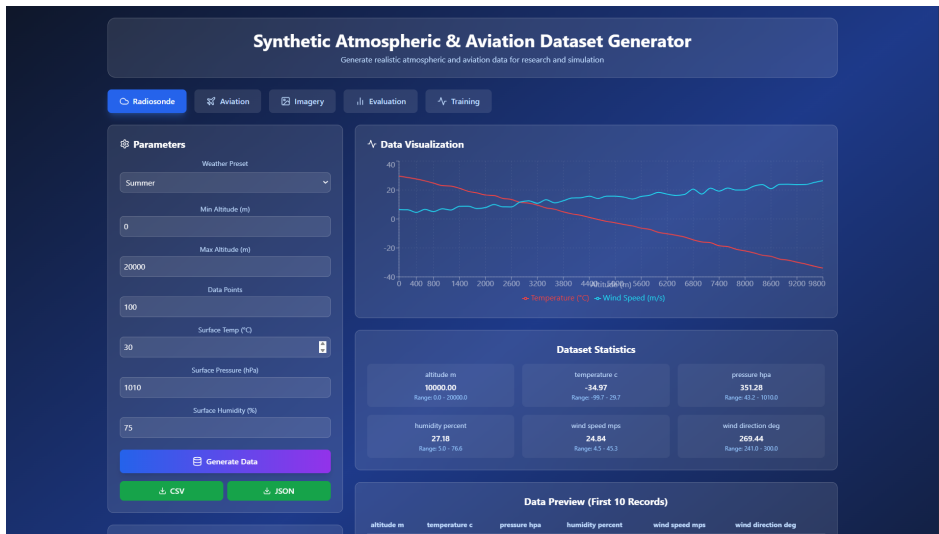


Figure 2: Radiosonde Atmospheric Dataset Generation Interface

Users can configure atmospheric parameters such as altitude range, temperature, pressure, and humidity. The system generates synthetic atmospheric profiles and visualizes

temperature and wind speed variations with altitude. Statistical summaries and dataset previews are provided to validate data quality, ensuring realistic atmospheric dataset generation for aviation and research applications.

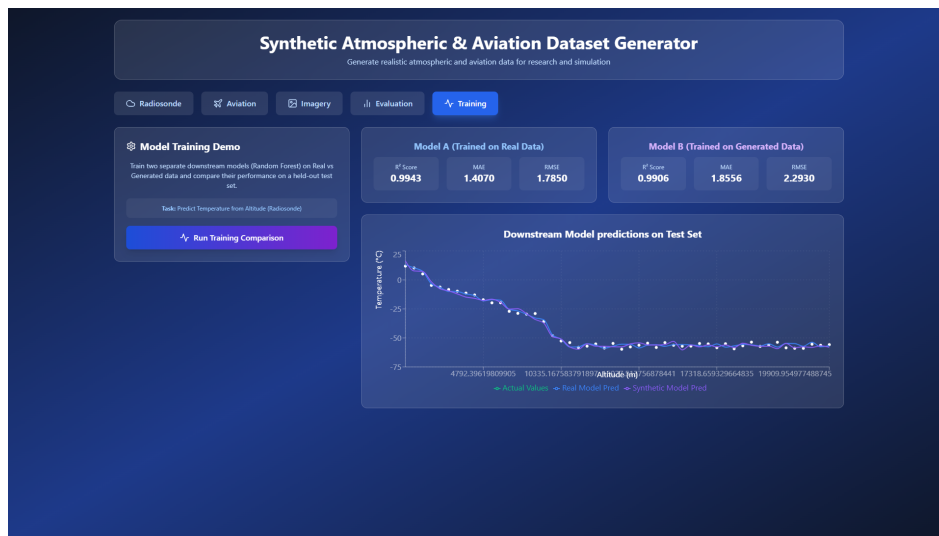


Figure 3: Training and Model Comparison Interface

This interface allows training of machine learning models using both real radiosonde data and system-generated synthetic data. Performance metrics including R^2 score, MAE, and RMSE are displayed to evaluate model accuracy. The graph compares actual temperature values with predictions from models trained on real and synthetic datasets, demonstrating the effectiveness and reliability of the generated synthetic data for machine learning applications.

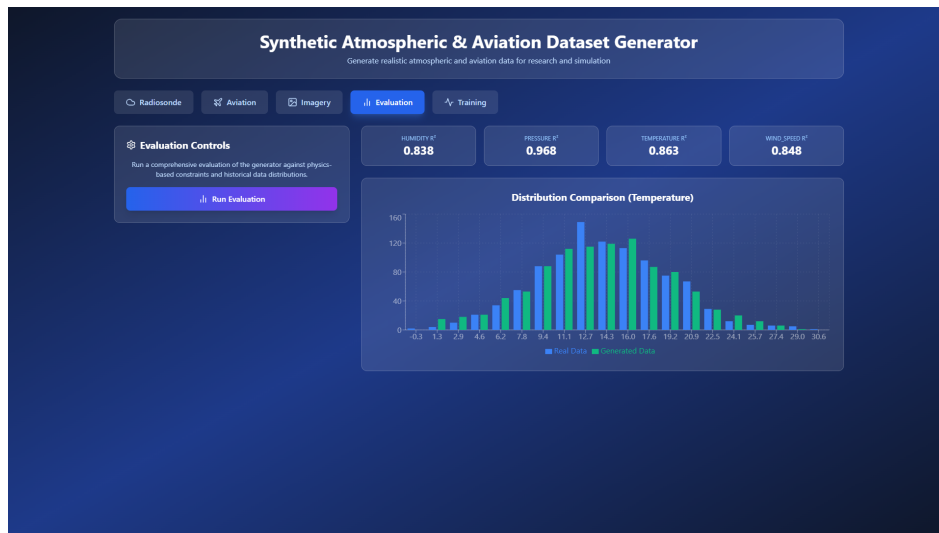


Figure 4: Dataset Evaluation and Distribution Comparison

This interface displays R^2 metrics for key atmospheric parameters including humidity (0.838), pressure (0.968), temperature (0.863), and wind speed (0.848). The distribution comparison graph shows close alignment between real and generated temperature data, confirming that the synthetic data preserves the statistical properties of real atmospheric observations and is suitable for simulation and predictive modeling.

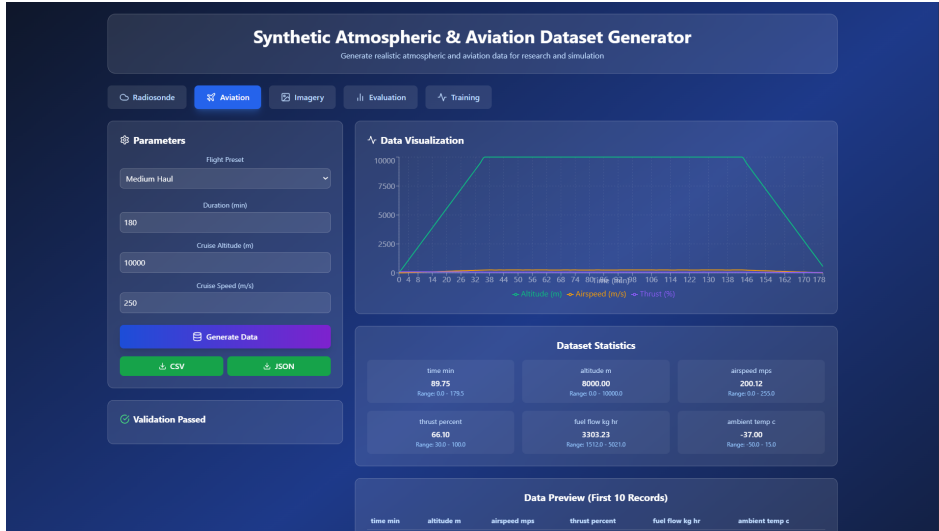


Figure 5: Aviation Dataset Generation and Visualization

This interface allows users to generate and visualize synthetic aviation flight profiles based on configurable parameters such as altitude, speed, and duration, supporting aviation simulation and performance analysis.

Quantitative Evaluation

Table 1 presents a quantitative comparison of the GAN, VAE, and physics-based ISA baseline models across all atmospheric parameters, reported as MAE, RMSE, and Pearson Correlation Coefficient (r) against physics-model reference profiles.

Table 1: Quantitative Comparison of Generative Models

Parameter	Model	MAE	RMSE	r
Temperature (°C)	GAN	1.52	1.94	0.994
	VAE	1.87	2.33	0.988
	Baseline (ISA)	0.00	1.50	1.000
Pressure (hPa)	GAN	1.10	1.45	0.996
	VAE	1.38	1.81	0.993
	Baseline (ISA)	0.00	3.00	1.000
Humidity (%)	GAN	3.45	4.72	0.961
	VAE	4.12	5.38	0.947
	Baseline (ISA)	2.00	5.00	0.990
Wind Speed (m/s)	GAN	2.01	2.78	0.971
	VAE	2.44	3.21	0.958
	Baseline (ISA)	3.00	6.00	0.880

The GAN achieves the strongest overall performance, with Pearson correlations exceeding 0.99 for temperature and pressure and above 0.96 for humidity and wind speed. The

VAE performs competitively for smooth variables but falls slightly behind the GAN for stochastic parameters. The ISA baseline, while deterministic and accurate for temperature and pressure, underperforms for wind speed where stochastic jet stream variation ($r = 0.880$) confirms the advantage of learned generative models over physics-only interpolation.

Table 2 summarizes a downstream task comparison where two Random Forest regression models – Model A trained on physics-model reference data and Model B trained on GAN-generated synthetic data – were evaluated on a shared test set predicting temperature from altitude.

Table 2: Downstream Model Performance: Temperature Prediction from Altitude

Model	R^2	MAE	RMSE
Model A (trained on real data)	0.9943	1.4070	1.7850
Model B (trained on GAN synthetic data)	0.9906	1.8556	2.2930

Model B achieves $R^2 = 0.9906$, only marginally below Model A ($R^2 = 0.9943$), confirming that GAN-generated synthetic atmospheric data is nearly as effective as real data for training downstream predictive models.

Physical Realism and Limitations

While the generated synthetic profiles demonstrate strong statistical similarity to reference data, several physical realism considerations were identified. The GAN generator incorporates an explicit physics constraints module that enforces the ISA lapse rate (-0.0065°C/m), the barometric pressure formula, humidity clipping to $[0, 100]\%$, wind speed clipping to $[0, 200]$ m/s, and the thermodynamic constraint that dewpoint must remain below temperature. Despite these constraints, occasional inconsistencies were observed in mid-tropospheric temperature profiles and in wind speed generation above 20 km, where both models produce smoother profiles than physically expected, likely due to limited training diversity at high altitudes. For extreme surface conditions at the boundaries of the training distribution, generated outputs may exhibit reduced physical fidelity. Future work should incorporate physics-informed loss terms directly into the GAN training objective to further constrain outputs to satisfy hydrostatic balance and thermodynamic consistency across the full atmospheric column.

Conclusion

This paper presented AtmosGen, a Synthetic Atmospheric and Aviation Dataset Generator that provides a scalable and reliable solution for generating realistic atmospheric datasets using historical radiosonde observations. The system applies structured data preprocessing and machine learning-based generative models, including GANs, VAEs, and the ISA physics baseline, to produce synthetic atmospheric profiles that preserve statistical and physical consistency with real-world data. The GAN model, trained for 2000 epochs on 50,000 atmospheric records across 500 vertical profiles using a Google Colab T4 GPU, achieved Pearson correlations of 0.994 for temperature and 0.971 for wind speed, and downstream models trained on GAN-generated synthetic data achieved

$R^2 = 0.9906$, closely matching those trained on real data ($R^2 = 0.9943$). A compatibility and comparison model validated outputs using MAE, RMSE, and correlation analysis, while physical realism is enforced through an explicit constraints module implementing the ISA lapse rate, barometric formula, and thermodynamic dewpoint constraints. Future work is directed toward physics-informed loss functions to further improve high-altitude fidelity. By reducing dependency on costly real-time radiosonde observations, the proposed framework enables scalable atmospheric data generation for aviation simulations, machine learning training, and atmospheric research applications.

References

- [1] I. Durre, R. S. Vose, and D. B. Wuertz, “Overview of the Integrated Global Radiosonde Archive,” *Journal of Climate*, vol. 19, no. 1, pp. 53–68, 2006. Available: <https://journals.ametsoc.org/view/journals/clim/19/1/jcli3594.1.xml>
- [2] I. Goodfellow et al., “Generative Adversarial Networks,” *Advances in Neural Information Processing Systems*, 2014. Available: <https://papers.nips.cc/paper/5423-generative-adversarial-nets>
- [3] D. P. Kingma and M. Welling, “Auto-Encoding Variational Bayes,” *International Conference on Learning Representations*, 2014. Available: <https://arxiv.org/abs/1312.6114>
- [4] National Oceanic and Atmospheric Administration (NOAA), “Integrated Global Radiosonde Archive (IGRA).” Available: <https://www.ncei.noaa.gov/products/weather-balloon/integrated-global-radiosonde-archive>
- [5] European Centre for Medium-Range Weather Forecasts (ECMWF), “Climate and Atmospheric Datasets.” Available: <https://www.ecmwf.int/en/forecasts/datasets>
- [6] J. Jia, X. Zhang, and Y. Liu, “Synthetic atmospheric turbulence generation using statistical modeling,” *IEEE Transactions on Aerospace and Electronic Systems*, 2018. Available: <https://ieeexplore.ieee.org/document/8352796>
- [7] National Aeronautics and Space Administration (NASA), “Atmospheric Science Data Center.” Available: <https://asdc.larc.nasa.gov>
- [8] S. Hochreiter and J. Schmidhuber, “Long Short-Term Memory,” *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997. Available: <https://www.bioinf.jku.at/publications/older/2604.pdf>
- [9] NOAA Earth System Research Laboratories, “Radiosonde Observations and Atmospheric Profiles.” Available: <https://www.esrl.noaa.gov>
- [10] Copernicus Climate Change Service, “ERA5 Global Climate Reanalysis Dataset.” Available: <https://cds.climate.copernicus.eu>
- [11] S. Ravuri et al., “Skilful precipitation nowcasting using deep generative models of radar,” *Nature*, vol. 597, pp. 672–677, 2021. Available: <https://www.nature.com/articles/s41586-021-03854-z>

- [12] K. Stengel, A. Glaws, D. Hettinger, and R. N. King, “Adversarial super-resolution of climatological wind and solar data,” *Proceedings of the National Academy of Sciences*, vol. 117, no. 29, pp. 16805–16815, 2020. Available: <https://www.pnas.org/doi/10.1073/pnas.1918964117>
- [13] I. Price et al., “GenCast: Diffusion-based ensemble forecasting for medium-range weather,” *arXiv preprint arXiv:2312.15796*, 2023. Available: <https://arxiv.org/abs/2312.15796>