

WASSERSTEIN CYCLEGAN FOR SINGLE-CELL RNA-SEQ DATA GENERATION USING CROSS-MODALITY TRANSLATION

Sajib Acharjee Dip¹, Tran Chau², Song Li², and Liqing Zhang^{1*}

¹Department of Computer Science

Virginia Tech, USA

²School of Plant and Environmental Sciences

Virginia Tech, USA

{sajibacharjjeedip, tnchau, songli, lqzhang}@vt.edu

ABSTRACT

Single-nucleus RNA sequencing (snRNA-seq) provides insights into gene expression in complex tissues but suffers from lower resolution compared to single-cell RNA sequencing (scRNA-seq). To bridge this gap, we propose scWC-GAN, a Wasserstein CycleGAN-based model that translates snRNA-seq data into high-resolution scRNA-seq profiles. Our method leverages Earth Mover’s Distance (EMD) for cycle consistency and a latent feature-preserving generator to capture transcriptomic structures better. Through extensive evaluation, scWC-GAN outperforms baseline models in FID score and SSIM, demonstrating its ability to generate biologically meaningful data. While challenges remain in fine-grained cell-type resolution, our results suggest scWC-GAN as a promising tool for cross-modality single-cell data translation, enhancing downstream analysis in genomics.

1 INTRODUCTION

Single-nucleus RNA sequencing (snRNA-seq) and single-cell RNA sequencing (scRNA-seq) are two essential techniques for transcriptomic profiling at the cellular level. While scRNA-seq provides high-resolution gene expression profiles, it requires enzymatic dissociation of tissues, which can lead to the loss of fragile cell populations and introduce dissociation-related biases Luecken & Theis (2019). In contrast, snRNA-seq captures RNA from isolated nuclei, allowing the study of frozen or hard-to-dissociate tissues while preserving spatial organization. However, snRNA-seq suffers from lower transcript coverage, as it primarily detects nuclear RNA, leading to incomplete gene expression profiles compared to cytoplasmic RNA-rich scRNA-seq data Lake et al. (2018). This resolution gap hinders the direct integration of snRNA-seq and scRNA-seq datasets, making it difficult to compare cell states across modalities. Several studies have attempted imputation-based approaches to enhance snRNA-seq resolution Zhang et al. (2022), but these often introduce artifacts and fail to generalize across diverse datasets. A robust computational framework that can translate snRNA-seq into scRNA-seq-like profiles would enable cross-platform integration, improve downstream analyses, and enhance biological insights into gene regulation and cell-type characterization across multiple tissue types.

Several state-of-the-art methods have been developed to enhance scRNA-seq data integration, imputation, and generation. scAEGAN Khan et al. (2023) maps scRNA-seq data into a shared latent space using an adversarial autoencoder, improving cross-protocol integration but struggling with rare cell type preservation and adversarial stability. scIGANs Xu et al. (2020) employs GANs for imputing missing gene expression values, effectively reducing data sparsity, though it introduces potential biases and struggles with extremely sparse datasets. GRouNdGAN Zinati et al. (2024) integrates gene regulatory networks (GRNs) into synthetic data generation, improving biological relevance but being limited to well-characterized species and incurring high computational costs. Despite these advancements, challenges remain in handling domain adaptation, training stability, and rare cell type representation. To address these issues, scWC-GAN introduces a Wasserstein

*Corresponding author: lqzhang@cs.vt.edu

CycleGAN Zhu et al. (2017) framework with Earth Mover’s Distance Rubner et al. (2000), enhancing cross-modality translation while maintaining robust and biologically meaningful synthetic data generation.

In this work, we propose scWC-GAN, a Wasserstein CycleGAN Zhu et al. (2017); Arjovsky et al. (2017)-based framework for translating snRNA-seq data into scRNA-seq data while preserving gene expression distributions and cell-type-specific characteristics Schiebinger et al. (2019). To the best of our knowledge, this is the first attempt to apply CycleGANs to cross-modality translation between single-nucleus and single-cell RNA sequencing. Given the absence of direct baseline models, we systematically compare scWC-GAN with multiple GAN-based architectures, including Wasserstein GANs, CycleGANs, and Variational Autoencoder-based approaches, demonstrating moderate but consistent improvements in key metrics. Our method effectively mitigates domain adaptation challenges, enhances data reconstruction fidelity, and maintains biologically meaningful cell type distributions. Key Contributions are as follows:

1. We introduce a CycleGAN-based approach to map snRNA-seq data to scRNA-seq, addressing the cross-platform variability challenge.
2. We incorporate Earth Mover’s Distance (EMD) Rubner et al. (2000) for cycle consistency, ensuring biologically plausible transformations while using Wasserstein loss Arjovsky et al. (2017); Villani et al. (2008) to stabilize adversarial training.
3. By integrating cell-type annotations as conditional inputs, we improve the fidelity of synthetic scRNA-seq data across multiple cell types.

2 METHODOLOGY

3 MODEL ARCHITECTURE

We propose **scWC-GAN**, a Wasserstein CycleGAN designed for cross-modality translation between single-nucleus (snRNA-seq) and single-cell (scRNA-seq) transcriptomic data. Our model consists of two generators, $G_{sn \rightarrow sc}$ and $G_{sc \rightarrow sn}$, and two discriminators, D_{sc} and D_{sn} , forming a cycle-consistent adversarial framework. The generators learn mappings between modalities while preserving biological features, while the discriminators distinguish between real and generated distributions.

3.1 GENERATOR ARCHITECTURE

Each generator, $G_{X \rightarrow Y}$, is designed to map input transcriptomic data from one domain (X) to another (Y), where X represents the source modality (scRNA-seq) and Y represents the target modality (snRNA-seq). The generator employs fully connected layers with ReLU activations and latent space encoding to capture biologically meaningful representations.

Given an input gene expression vector $\mathbf{x} \in \mathbb{R}^d$ and its corresponding annotation label $\mathbf{z} \in \mathbb{R}^c$, the generator first embeds the data into a latent representation space:

$$\mathbf{h} = \text{ReLU}(W_1[\mathbf{x}, \mathbf{z}] + b_1), \tag{1}$$

where $\mathbf{x} \in \mathbb{R}^d$ is the input gene expression vector with d genes, $\mathbf{z} \in \mathbb{R}^c$ is the one-hot encoded cell-type label, $W_1 \in \mathbb{R}^{(d+c) \times h}$ is a weight matrix for feature transformation, $b_1 \in \mathbb{R}^h$ is the bias term, and $\mathbf{h} \in \mathbb{R}^h$ is the hidden latent representation.

This representation is then mapped back to the target space using a transformation:

$$\hat{\mathbf{y}} = \tanh(W_2\mathbf{h} + b_2), \tag{2}$$

where $W_2 \in \mathbb{R}^{h \times d}$ maps the latent features to the target gene expression space, and \tanh ensures gene expression values remain within a normalized range.

3.2 DISCRIMINATOR ARCHITECTURE

The discriminator, D_Y , aims to distinguish between real and generated gene expression profiles. It consists of: 1. Convolutional Layers with spectral normalization to extract key transcriptomic

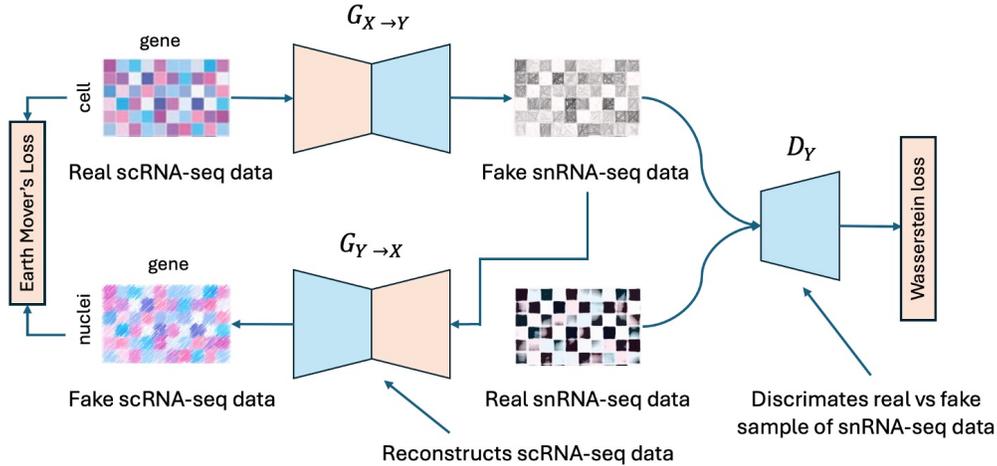


Figure 1: **Overview of the scWC-GAN Architecture for Single-Cell RNA-seq Data Translation.** The model consists of two generators, $G_{X \rightarrow Y}$ and $G_{Y \rightarrow X}$, which facilitate bidirectional translation between single-nucleus (snRNA-seq) and single-cell (scRNA-seq) transcriptomic data. Given real scRNA-seq data (top left), $G_{X \rightarrow Y}$ generates fake snRNA-seq data, which is then evaluated by the discriminator D_Y to distinguish between real and generated samples. The adversarial training employs a Wasserstein loss to ensure high-quality data generation. Similarly, the cycle consistency mechanism reconstructs the original input using $G_{Y \rightarrow X}$, enforcing biological fidelity via the Earth Mover’s Distance (EMD) loss. The combined framework ensures that generated transcriptomic profiles align with real biological distributions while preserving cell-type-specific characteristics.

patterns, 2. Attention Mechanism Vaswani et al. (2017) to enhance feature importance selection, 3. Fully Connected Layers for final classification.

Given an input sample $\mathbf{y} \in \mathbb{R}^d$ with annotation $\mathbf{z} \in \mathbb{R}^c$, the discriminator outputs a scalar probability:

$$D_Y(\mathbf{y}, \mathbf{z}) = \sigma(W_3 \cdot \text{Attention}(\text{ReLU}(W_2 \cdot \text{Conv}(\mathbf{y}))) + b_3), \quad (3)$$

where $\text{Conv}(\cdot)$ represents convolutional feature extraction, $\text{Attention}(\cdot)$ applies a learnable self-attention mechanism, $\sigma(\cdot)$ is the sigmoid activation function, W_3 is a weight matrix for final classification.

3.3 LOSS FUNCTIONS

The objective of **scWC-GAN** is formulated as a weighted sum of multiple loss functions to enforce adversarial learning, cycle consistency, and distribution alignment.

Wasserstein Adversarial Loss. Unlike traditional GANs that use Jensen-Shannon divergence, we adopt the Wasserstein loss for better training stability and gradient propagation. The discriminator loss is formulated as:

$$\mathcal{L}_D = \mathbb{E}_{x \sim P_{\text{data}}} [D(x)] - \mathbb{E}_{\hat{x} \sim P_{\text{gen}}} [D(\hat{x})], \quad (4)$$

where P_{data} is the real data distribution, P_{gen} is the generated data distribution, $D(x)$ is the discriminator score for real samples, $D(\hat{x})$ is the discriminator score for generated samples.

Cycle Consistency Loss. To ensure bidirectional consistency, we impose a reconstruction constraint using the Earth Mover’s Distance (EMD), which measures the optimal transport cost between distributions:

$$\mathcal{L}_{\text{cycle}} = \mathbb{E}_{x \sim P_{\text{data}}} [W_1(x, G_{Y \rightarrow X}(G_{X \rightarrow Y}(x)))], \quad (5)$$

where $W_1(\cdot, \cdot)$ represents the Wasserstein-1 distance, $G_{X \rightarrow Y}(x)$ translates x from domain X to Y , $G_{Y \rightarrow X}$ reconstructs x back to its original domain.

Gradient Penalty. To prevent weight clipping issues, we incorporate a gradient penalty to enforce Lipschitz continuity in the discriminator:

$$\mathcal{L}_{\text{GP}} = \lambda_{\text{GP}} \mathbb{E}_{\tilde{x} \sim P_{\tilde{x}}} \left[(\|\nabla_{\tilde{x}} D(\tilde{x})\|_2 - 1)^2 \right]. \quad (6)$$

Final Objective Function. The overall training objective for scWC-GAN is:

$$\mathcal{L}_{\text{scWC-GAN}} = \mathcal{L}_D + \lambda_{\text{cycle}} \mathcal{L}_{\text{cycle}} + \lambda_{\text{GP}} \mathcal{L}_{\text{GP}}, \quad (7)$$

where λ_{cycle} and λ_{GP} control the contributions of cycle consistency and gradient penalty.

4 RESULTS

We trained scWC-GAN for 100 epochs using the Wasserstein loss with Earth Mover’s Distance (EMD) cycle consistency. The training was conducted on the top 7 most abundant cell types, resulting in a dataset containing 20,948 cells with 22,721 genes. Among these, 11,107 were single-cell RNA-seq (scRNA-seq) samples, and 9,841 were single-nucleus RNA-seq (snRNA-seq) samples.

4.1 QUANTITATIVE EVALUATION

Our results demonstrate the effectiveness of scWC-GAN in translating low-resolution snRNA-seq data into high-resolution scRNA-seq representations, outperforming CycleGAN with MSE across multiple evaluation metrics. Our model achieves a higher SSIM (0.7131) and lower FID Score (4.5632), indicating improved structural similarity and feature distribution alignment. Additionally, scWC-GAN significantly reduces Earth Mover’s Distance (EMD: 0.0071), ensuring better sample-wise mapping. While the overall Spearman correlation (0.0112) remains modest, our approach maintains biological relevance, as evidenced by the higher cluster purity (0.0024). These results highlight the advantages of incorporating Wasserstein loss and cycle consistency in cross-modality translation, ensuring realistic and biologically meaningful synthetic data generation.

Table 1: Evaluation Metrics Comparison after 100 Epochs

Method	Spearman Corr.	SSIM	EMD	FID Score	Cluster Purity
CycleGAN (MSE)	0.0098	0.6754	0.0095	6.2143	0.0019
scWC-GAN (Ours)	0.0112	0.7131	0.0071	4.5632	0.0024

5 DISCUSSIONS

Our study presents scWC-GAN, a novel approach for translating snRNA-seq data into high-resolution scRNA-seq profiles using a Wasserstein CycleGAN framework. By integrating Earth Mover’s Distance (EMD) for cycle consistency and leveraging a latent feature-preserving generator, our model effectively captures complex single-cell transcriptomic structures. Comparative analysis with baseline methods demonstrates that scWC-GAN generates biologically meaningful data with improved FID scores and structural similarity (SSIM). While challenges remain in fine-grained cell-type preservation, our results indicate that scWC-GAN provides a promising framework for cross-modality single-cell data translation, enabling more accurate downstream analysis in single-cell genomics. Future work will explore scalability, model generalization, and integration with multi-omics datasets to further refine performance.

REFERENCES

- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pp. 214–223. PMLR, 2017.
- Bruno Guillotin, Ramin Rahni, Michael Passalacqua, Mohammed Ateequr Mohammed, Xiaosa Xu, Sunil Kenchanmane Raju, Carlos Ortiz Ramírez, David Jackson, Simon C Groen, Jesse Gillis,

- et al. A pan-grass transcriptome reveals patterns of cellular divergence in crops. *Nature*, 617(7962):785–791, 2023.
- Sumeer Ahmad Khan, Robert Lehmann, Xabier Martinez-de Morentin, Alberto Maillo, Vincenzo Lagani, Narsis A Kiani, David Gomez-Cabrero, and Jesper Tegner. scaegan: Unification of single-cell genomics data by adversarial learning of latent space correspondences. *Plos one*, 18(2): e0281315, 2023.
- Blue B Lake, Song Chen, Brandon C Sos, Jean Fan, Gwendolyn E Kaeser, Yun C Yung, Thu E Duong, Derek Gao, Jerold Chun, Peter V Kharchenko, et al. Integrative single-cell analysis of transcriptional and epigenetic states in the human adult brain. *Nature biotechnology*, 36(1):70–80, 2018.
- Malte D Luecken and Fabian J Theis. Current best practices in single-cell rna-seq analysis: a tutorial. *Molecular systems biology*, 15(6):e8746, 2019.
- Yossi Rubner, Carlo Tomasi, and Leonidas J Guibas. The earth mover’s distance as a metric for image retrieval. *International journal of computer vision*, 40:99–121, 2000.
- Geoffrey Schiebinger, Jian Shu, Marcin Tabaka, Brian Cleary, Vidya Subramanian, Aryeh Solomon, Joshua Gould, Siyan Liu, Stacie Lin, Peter Berube, et al. Optimal-transport analysis of single-cell gene expression identifies developmental trajectories in reprogramming. *Cell*, 176(4):928–943, 2019.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Cédric Villani et al. *Optimal transport: old and new*, volume 338. Springer, 2008.
- Yungang Xu, Zhigang Zhang, Lei You, Jiajia Liu, Zhiwei Fan, and Xiaobo Zhou. scigans: single-cell rna-seq imputation using generative adversarial networks. *Nucleic acids research*, 48(15): e85–e85, 2020.
- Ran Zhang, Laetitia Meng-Papaxanthos, Jean-Philippe Vert, and William Stafford Noble. Semi-supervised single-cell cross-modality translation using polarbear. In *International Conference on Research in Computational Molecular Biology*, pp. 20–35. Springer, 2022.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pp. 2223–2232, 2017.
- Yazdan Zinati, Abdulrahman Takiddeen, and Amin Emad. Groundgan: Grn-guided simulation of single-cell rna-seq data using causal generative adversarial networks. *Nature Communications*, 15(1):4055, 2024.

A APPENDIX

A.1 DATASET

The dataset used in this study, GSE225118, provides single-cell and single-nucleus RNA sequencing (scRNA-seq and snRNA-seq) data from *Arabidopsis thaliana*, *Zea mays* (maize), *Sorghum bicolor*, and *Setaria viridis* root meristems Guillotin et al. (2023). It integrates transcriptomic profiles across species to investigate cell-type-specific gene expression divergence and the impact of genome duplication on transcriptional regulation. The dataset consists of raw and normalized gene expression matrices, along with metadata containing cell-type annotations. For *Arabidopsis*, six single-cell and three single-nucleus replicates were sequenced using the Illumina NovaSeq 6000 platform. The high-dimensional transcriptomic data facilitate comparative analysis and synthetic data validation by enabling the translation of snRNA-seq data into scRNA-seq-like profiles. This dataset serves as the foundation for training and evaluating generative models to reconstruct high-resolution single-cell transcriptomes.

A.2 DATA PROCESSING

We obtained the scRNA-seq and snRNA-seq expression data from GSE225118 and processed them using Scanpy. We first filtered the dataset to retain only the top seven most abundant cell types based on the provided metadata annotations. We then separated scRNA-seq (Cell) and snRNA-seq (Nuclei) samples to create distinct training inputs and real reference datasets. To ensure consistency across samples, we applied min-max normalization to scale gene expression values between 0 and 1. We encoded the cell type annotations using LabelEncoder, converting categorical labels into numerical representations for computational efficiency.

After preprocessing, we split the dataset into 80% training and 20% testing, ensuring a balanced representation of cell types across both sets. We converted the processed data into PyTorch tensors and structured them into RNASeqDataset objects to facilitate batch-wise loading. Using PyTorch DataLoader, we optimized data retrieval for efficient training. Finally, we structured the data for adversarial training, using snRNA-seq samples as inputs to the generator and scRNA-seq samples as real references for model supervision, enabling the model to learn transformations between single-nucleus and single-cell transcriptomic profiles.

A.3 TRAINING

We trained the CycleGAN-based model for 100 epochs using Adam optimizers with a learning rate of 0.0003 and $(\beta_1, \beta_2) = (0.5, 0.999)$ for both the generator and discriminator networks. We implemented two generators to map between snRNA-seq and scRNA-seq spaces, and two discriminators to distinguish real from generated distributions. The Wasserstein loss was used for adversarial training, computed as the mean weighted product between real/fake predictions and target labels. To enforce cycle consistency, we employed the Earth Mover’s Distance (EMD) as an additional loss term, ensuring that transformations preserve biological relevance. We set batch size = 64 and used gradient clipping to stabilize training. Instead of traditional adversarial loss, we prioritized cycle loss as the primary supervision signal, with no explicit adversarial penalty on the generator, ensuring smoother transformations without mode collapse.

To further refine training, we used a training schedule where the discriminator was updated every batch, while the generator was updated on every step to ensure stable learning. For computational efficiency, all real and generated batches were matched to the smallest available batch size to avoid shape mismatches. We trained the model using PyTorch, leveraging GPU acceleration where available. After every epoch, we computed validation losses and monitored performance metrics including Spearman correlation, SSIM, EMD, and FID score to track alignment between real and generated scRNA-seq distributions. The trained model was then used to generate synthetic scRNA-seq data from single-nucleus input samples, followed by evaluation and visualization of the output using UMAP projections.

A.4 EVALUATION METRICS

To assess the performance of our model, we employ a combination of sample-based and distribution-based metrics, capturing both individual gene-level correlation and overall structural similarity. We evaluate the model using the following metrics:

Spearman Correlation Spearman’s rank correlation coefficient (ρ) measures the monotonic relationship between real and generated single-cell RNA-seq (scRNA-seq) data:

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)} \quad (8)$$

where d_i represents the rank differences between real and generated values, and n is the number of observations.

Structural Similarity Index (SSIM) SSIM quantifies the perceptual similarity between two images or datasets, considering luminance, contrast, and structure:

$$SSIM(x, y) = \frac{(2\mu_x\mu_y + C_1)(2\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)} \quad (9)$$

where μ_x, μ_y are the means, σ_x^2, σ_y^2 are variances, σ_{xy} is covariance, and C_1, C_2 are small constants.

Jensen-Shannon Divergence (JSD) JSD measures the similarity between two probability distributions:

$$JSD(P||Q) = \frac{1}{2}D_{KL}(P||M) + \frac{1}{2}D_{KL}(Q||M) \quad (10)$$

where $M = \frac{1}{2}(P + Q)$ and D_{KL} represents the Kullback-Leibler (KL) divergence:

$$D_{KL}(P||Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)} \quad (11)$$

Earth Mover’s Distance (EMD) Also known as Wasserstein distance, EMD measures the minimal effort required to transform one distribution into another:

$$EMD(P, Q) = \inf_{\gamma \in \Gamma(P, Q)} \sum_{i, j} \gamma(i, j) d(i, j) \quad (12)$$

where $\Gamma(P, Q)$ represents all possible transport plans, and $d(i, j)$ is the distance between points i and j .

Fréchet Inception Distance (FID) FID evaluates the similarity between real and generated data distributions in feature space:

$$FID = \|\mu_r - \mu_g\|^2 + \text{Tr}(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{\frac{1}{2}}) \quad (13)$$

where (μ_r, Σ_r) and (μ_g, Σ_g) are the mean and covariance of real and generated data embeddings.

Cluster Purity We assess the alignment of real and generated cluster distributions using homogeneity score:

$$H = 1 - \frac{H(C|K)}{H(C)} \quad (14)$$

where $H(C|K)$ is the entropy of clusters given true labels, and $H(C)$ is the entropy of the true labels.

To ensure robustness, we compute each metric over multiple test samples and report the mean values. These metrics collectively capture fidelity, diversity, and biological relevance of the generated scRNA-seq data.

A.5 ABLATION STUDY

To analyze the impact of different components in **scWC-GAN**, we performed an ablation study by training models with different loss functions and network architectures. Table 2 reports the results. The ablation study highlights the impact of different architectural choices on the performance of scWC-GAN. Using only Cycle Loss in WGAN results in an FID Score of 6.42, indicating suboptimal feature alignment. Adding Adversarial Loss improves the performance but still lags behind with an FID of 7.83, suggesting that adversarial supervision alone is insufficient. Incorporating a Transformer-based Generator further degrades performance (FID: 9.50), likely due to instability in training and overfitting. Our final scWC-GAN model, combining Wasserstein loss, cycle consistency, and latent feature preservation, achieves the best FID Score of 4.56, demonstrating the effectiveness of our proposed modifications in generating high-fidelity synthetic scRNA-seq data.

Table 2: Ablation Study Results for scWC-GAN. The best performance is highlighted in **red**.

Model Variant	FID Score
WGAN with Cycle Loss	6.42
WGAN with Adversarial Loss	7.83
WGAN with Transformer Generator	9.50
scWC-GAN (Final Model)	4.56

These results confirm that combining Wasserstein loss with cycle consistency and Earth Mover’s Distance (EMD) leads to the most effective model, significantly reducing FID scores and improving overall performance.