

## Article

# Physically Informed Synthetic Data Generation and U-Net Generative Adversarial Network for Palimpsest Reconstruction

Jose L. Salmeron <sup>1,\*</sup>  and Eva Fernandez-Palop <sup>2</sup><sup>1</sup> School of Engineering, CUNEF University, 28040 Madrid, Spain<sup>2</sup> Department of Arts, Universidad de Zaragoza, 44003 Teruel, Spain; e.fernandez@unizar.es

\* Correspondence: joseluis.salmeron@cunef.edu

## Abstract

This paper introduces a novel adversarial learning framework for reconstructing hidden layers in historical palimpsests. Recovering text hidden in historical palimpsests is complicated by various artifacts, such as ink diffusion, degradation of the writing substrate, and interference between overlapping layers. To address these challenges, the authors of this paper combine a synthetic data generator grounded in physical modeling with three generative architectures: a baseline VAE, an improved variant with stronger regularization, and a U-Net-based GAN that incorporates residual pathways and a mixed loss strategy. The synthetic data engine aims to emulate key degradation effects—such as ink bleeding, the irregularity of parchment fibers, and multispectral layer interactions—using stochastic approximations of underlying physical processes. The quantitative results suggest that the U-Net-based GAN architecture outperforms the VAE-based models by a notable margin, particularly in scenarios with heavy degradation or overlapping ink layers. By relying on synthetic training data, the proposed method facilitates the non-invasive recovery of lost text in culturally important documents, and does so without requiring costly or specialized imaging setups.

**Keywords:** palimpsest reconstruction; generative adversarial networks; deep learning; synthetic data generation; cultural heritage; multispectral imaging

**MSC:** 68T01; 68T05; 68T07; 68T30



Academic Editor: Raymond Lee

Received: 1 July 2025

Revised: 16 July 2025

Accepted: 17 July 2025

Published: 18 July 2025

**Citation:** Salmeron, J.L.; Fernandez-Palop, E. Physically Informed Synthetic Data Generation and U-Net Generative Adversarial Network for Palimpsest Reconstruction. *Mathematics* **2025**, *13*, 2304. <https://doi.org/10.3390/math13142304>

**Copyright:** © 2025 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

Palimpsests, manuscripts with erased underlying texts, represent invaluable cultural heritage artifacts [1]. These multilayered documents, where primary texts were eradicated through chemical, mechanical, or optical means to enable surface reuse, are of significant paleographic interest. These artifacts constitute complex material systems wherein original substrates (typically parchment or vellum) undergo deliberate degradation processes followed by secondary inscription, resulting in stratified textual interventions separated temporally and chemically.

Contemporary multispectral imaging and computational analysis techniques are used to attempt spectral separation of these superimposed textual layers, yet remain fundamentally constrained by nonlinear interactions between (1) residual chromophores from eradicated texts, (2) later-overwritten inscriptions, and (3) time-dependent physicochemical alterations of the writing medium. Quantitative analysis of these interactions—particularly regarding signal modulation through fiber occlusion, character fusion, and contrast

degradation—requires rigorous metric validation to address the inherent information loss in such re-mediated textual corpora.

Traditional multispectral imaging (MSI) techniques face limitations in terms of accessibility and resolution [2]. Recent computational approaches [3] have shown promise but struggle with complex degradation patterns.

Generative models offer potential solutions, yet their application to palimpsests remains underexplored. Although VAEs provide probabilistic frameworks [4], they often yield blurred reconstructions. GANs generate sharper outputs [5] but face training instability in low-data regimes. Our proposed method bridges this gap through three key contributions:

- A physically informed synthetic generator modeling parchment degradation.
- A novel GAN architecture with asymmetric skip connections.
- The first comparative analysis of generative models for palimpsest reconstruction.

The remainder of the paper is organized as follows. Section 2 analyzes related work. Section 3 details the methodological approach. The results are shown in Section 4. Finally, the conclusions are presented in Section 5.

## 2. Related Works

Traditional approaches to palimpsest reconstruction have primarily relied on specialized imaging techniques. Multispectral imaging (MSI) [2] remains the gold standard, capturing reflectance properties across spectral bands to enhance latent text visibility. More advanced methods include X-ray fluorescence (XRF) imaging [1] and optical coherence tomography [6], which enable non-invasive material characterization. While effective, these hardware-dependent approaches require expensive instrumentation and controlled environments, limiting accessibility for cultural heritage institutions. Computational enhancements like principal component analysis (PCA) and independent component analysis (ICA) have been applied to MSI data [3], but remain constrained by physical capture limitations.

A recent review by Perino et al. [7] highlights the progress made in the digital restoration of historical manuscripts, with a particular focus on recovering writings obscured by erasure, fading, carbonization, and the natural effects of aging. These challenges have historically rendered many texts illegible to the naked eye. However, contemporary technological advancements have markedly enhanced our capacity to recover and analyze this hidden written heritage, enabling access to cultural materials that had long been beyond the reach of modern scholarship.

To address the challenge of manuscript reconstruction, Jampour [3] applies an effective image inpainting technique based on a generative model. The proposed method leverages a Latent Diffusion Model (LDM) backbone, incorporating key modifications to the conditioning mechanism that allow the model to effectively utilize contextual information from the surrounding regions of the mask. To further enhance the generation process, the author provides an initial approximation of the masked region's pixels as a starting condition.

Recent advances have explored data-driven approaches to historical document analysis. Convolutional Neural Networks (CNNs) have been applied to tasks that include ink bleeding reduction [8] and document binarization [9]. Christlein et al. [3] pioneered CNN-based methods specifically for palimpsests, although their approach required paired real-world multispectral captures. Generative Adversarial Networks (GANs) have shown promise in related domains, with Isola et al. [10] demonstrating image-to-image translation capabilities using conditional adversarial networks. However, direct application to palimpsests faces challenges due to the scarcity of training data with verified ground truth.

Probabilistic approaches using variational autoencoders (VAEs) [4] have been explored for artifact reconstruction. Ronneberger's U-Net architecture [11] has been adapted for

manuscript fragment reassembly, while attention mechanisms [12] have improved feature localization in degraded regions. Self-attention GANs [13] have achieved state-of-the-art results in high-fidelity image generation, although their application to layered document reconstruction remains unexplored. The adversarial training framework by Heusel et al. [14] provides a theoretical foundation for the proposed optimization approach.

Starynska et al. [15] propose revealing under-text completely using deep generative networks, by leveraging the prior spatial information of the under-text script. To optimize the under-text, the authors mimic the process of palimpsest creation, generating the under-text from a separately trained generative network to match it to the palimpsest image after mixing it with foreground text.

To address data scarcity, researchers have developed synthetic degradation models. Chen et al. [8] simulated document aging through random noise and blur operations, while Mitra et al. [6] modeled spectral interactions in multisensor systems. This research advances this paradigm through physically informed degradation processes that explicitly model the following:

- Ink diffusion: parameterized by Fick's second law of diffusion [16].
- Parchment structure: biomechanical fiber modeling.
- Spectral superposition: wavelength-dependent layer interactions.

Recent advances in AI, particularly in machine learning and deep learning, have demonstrated significant potential in enhancing data assimilation (DA) and uncertainty quantification (UQ) for problems in Earth sciences. Two notable applications are emerging in the context of geologic carbon storage (GCS) and seismic velocity model building (VMB).

The first line of research addresses the optimization of history matching in GCS scenarios by integrating surrogate AI models into hybrid data assimilation frameworks. Specifically, surrogate-assisted hybrid methods such as SH-ESMDA (Ensemble Smoother with Multiple Data Assimilation) and SH-RML (Regularized Maximum Likelihood) have been proposed to reduce the computational burden of traditional simulations while maintaining robust inference performance [17]. These approaches employ deep neural operators such as Fourier Neural Operators (FNOs) and Transformer-based U-Nets (T-UNets) as fast approximators of complex subsurface flow models, significantly accelerating the assimilation process and improving match quality with historical data.

Another research efforts focuses on the generation of high-fidelity synthetic data for deep learning inversion methods in seismic VMB. This effort is motivated by the scarcity of labeled data and the need for generalizable models across geologically diverse scenarios. A comprehensive dataset generation workflow is proposed that synthesizes geophysical models containing varied geological patterns, including horizontal stratified layers, folded structures, and complex salt dome intrusions. These synthetic datasets are then used to train DL-based inversion networks capable of producing reliable subsurface reconstructions under realistic noise and acquisition conditions [18].

These contributions illustrate how AI-driven methodologies can substantially improve the efficiency and reliability of geophysical simulation and subsurface interpretation workflows. By leveraging domain-informed architectures and physically consistent training regimes, these AI surrogates not only reduce computational costs but also facilitate scalable uncertainty-aware modeling in highly complex geological environments [19,20].

Existing methods exhibit three key limitations: (1) dependence on specialized hardware, (2) the absence of generalized degradation models, and (3) the inadequate handling of overlapping script features. This paper bridges these gaps through a purely computational framework that combines physical simulation with adversarial learning, enabling reconstruction without multispectral inputs.

### 3. Methodological Approach

#### 3.1. Synthetic Data Generation Framework

To overcome the scarcity of paired palimpsest datasets, a synthetic data generator grounded in physical modeling was developed. A comparable methodology was adopted by Zheng et al. in [21]. The generator simulates manuscript degradation by modeling stochastic physical processes and integrates three key components (parchment texture modeling, ink degradation, and multispectral layer superposition). The pseudo-code for physically informed palimpsest sample generation is shown at Algorithm 1.

---

**Algorithm 1:** Physically Informed Palimpsest Sample Generation
 

---

**Input** : Image size  $S = (H, W)$ , script list  $\mathcal{S}$ , font map  $\mathcal{F}$ , corpus  $\mathcal{C}$ , degradation parameters  $\theta$ ;

**Output**: RGB image with two text layers and degradation: combined, underlying, overwritten

- 1 Randomly sample two distinct scripts:  $s_1, s_2 \leftarrow \text{sample}(\mathcal{S}, 2)$ ;
  - 2 Generate random text:  $t_1 \leftarrow \text{sample\_text}(\mathcal{C}[s_1])$ ,  $t_2 \leftarrow \text{sample\_text}(\mathcal{C}[s_2])$ ;
  - 3 Render first text layer:  $L_1 \leftarrow \text{render\_text}(t_1, s_1, \mathcal{F}[s_1])$ ;
  - 4 Apply physical degradation:  $L_1 \leftarrow \text{degrade}(L_1, \theta)$ ;
  - 5 Render second text layer:  $L_2 \leftarrow \text{render\_text}(t_2, s_2, \mathcal{F}[s_2])$ ;
  - 6 Apply physical degradation:  $L_2 \leftarrow \text{degrade}(L_2, \theta)$ ;
  - 7 Generate realistic parchment texture:  $P \leftarrow \text{generate\_parchment}(S)$ ;
  - 8 Randomly sample transparency coefficient  $\alpha \in [0.3, 0.6]$ ;
  - 9 Blend overwritten text:  $B \leftarrow P \cdot (1 - L_2) + 0.2 \cdot L_2$ ;
  - 10 Blend underlying text:  $B \leftarrow B \cdot (1 - \alpha \cdot L_1) + \alpha \cdot L_1$ ;
  - 11 Add Gaussian noise:  $B \leftarrow \text{clip}(B + \mathcal{N}(0, 0.02), 0, 1)$ ;
  - 12 Enhance contrast of text layers:  $L_1 \leftarrow \text{clip}(L_1 \cdot 3.0, 0, 1)$ ,  $L_2 \leftarrow \text{clip}(L_2 \cdot 2.5, 0, 1)$ ;
  - 13 Return combined =  $B$ , underlying =  $L_1$ , overwritten =  $L_2$ , and scripts =  $(s_1, s_2)$ ;
- 

For parchment texture modeling, the authors simulate parchment fibers as anisotropic structures using Equation (1).

$$T(x, y) = \underbrace{\mathcal{U}(0.7, 0.9)}_{\text{base intensity}} \otimes \underbrace{\mathcal{F}_{(l, \theta)}}_{\text{fiber model}} + \underbrace{\mathcal{S}_\lambda \otimes \mathcal{U}(0.3, 0.7)}_{\text{stain model}} \quad (1)$$

where  $T(x, y)$  denotes the simulated parchment texture at position  $(x, y)$ , and the symbol  $\mathcal{U}(a, b)$  represents a uniform distribution in the range  $[a, b]$ . The intensity ranges  $[0.7, 0.9]$  and  $[0.3, 0.7]$  in Equation (1) are selected to replicate the visual characteristics of real historical parchment surfaces. The base intensity, drawn from the uniform distribution  $\mathcal{U}(0.7, 0.9)$ , models the general reflectance of the parchment background. This range reflects the typical brightness of parchment, which is light in tone but not uniformly white, capturing subtle variations caused by aging and material texture. Moreover, the operator  $\otimes$  indicates a spatial modulation (i.e., pixel-wise multiplication) between two image components. The fiber model  $\mathcal{F}_{(l, \theta)}$  produces anisotropic linear structures, with fiber lengths  $l$  drawn from a uniform distribution  $\mathcal{U}(10, 30)$  pixels, and orientation angles  $\theta$  from  $\mathcal{U}(0, \pi)$ . The stain model  $\mathcal{S}_\lambda$  introduces spatially sparse degradations, modeled using a Poisson distribution with the parameter  $\lambda = 0.02$ . The stain intensity is sampled from  $\mathcal{U}(0.3, 0.7)$  to represent darker regions associated with surface imperfections, such as ink bleed-through, wear, or biological degradation. This range ensures sufficient contrast with the base while avoiding unrealistically dark

artifacts. Together, these empirically chosen intervals approximate the statistical reflectance properties observed in digitized historical manuscripts [1], and they can be adjusted to simulate different parchment conditions or lighting environments. This formulation captures the non-uniform reflectance and irregularities typical of historical parchment surfaces [1].

Ink degradation is modeled with a coupled diffusion–bleeding approach, as shown in Equation (2). Ink degradation is modeled as a combination of diffusion, bleeding, and additive noise. This is expressed in Equation (2) as follows:

$$I'(x, y) = \underbrace{I(x, y) * G_\sigma}_{\text{diffusion}} + \underbrace{\beta I(x, y)}_{\text{bleeding}} + \underbrace{\mathcal{N}(0, \sigma_n)}_{\text{noise}} \quad (2)$$

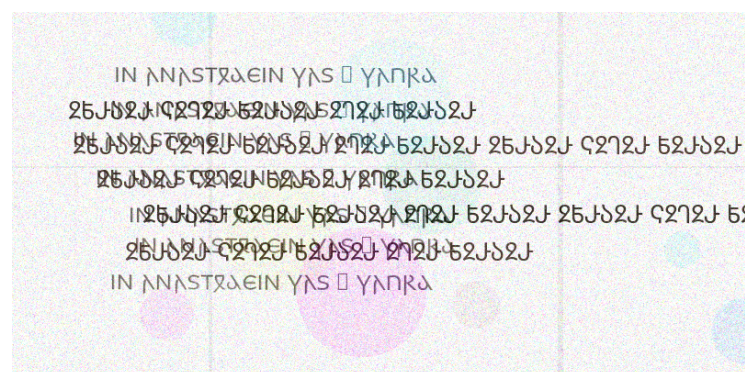
where  $I(x, y)$  denotes the original ink intensity at pixel  $(x, y)$ , and  $I'(x, y)$  is the degraded ink after the simulation. The symbol  $*$  indicates 2D convolution. The term  $G_\sigma$  represents a Gaussian kernel with a standard deviation of  $\sigma = 1.2$  and a size of  $5 \times 5$ , used to simulate isotropic diffusion. Convolution with a Gaussian kernel approximates Fickian diffusion, since the solution to the diffusion equation under homogeneous, isotropic conditions is a Gaussian function. The second term models ink bleeding, scaled by a coefficient  $\beta = 0.05$ , and the third term adds Gaussian noise with zero mean and a standard deviation of  $\sigma_n = 0.02$ . This formulation follows prior approaches to simulate ink degradation effects observed in historical documents [6].

For multispectral layer superposition, the final palimpsest combines layers through the Equation (3).

$$P_{final} = \underbrace{\alpha L_{under}}_{\text{subtext}} + \underbrace{(1 - \alpha)(L_{over} \otimes P_{parchment})}_{\text{overtext}} + \mathcal{N}(0, 0.02) \quad (3)$$

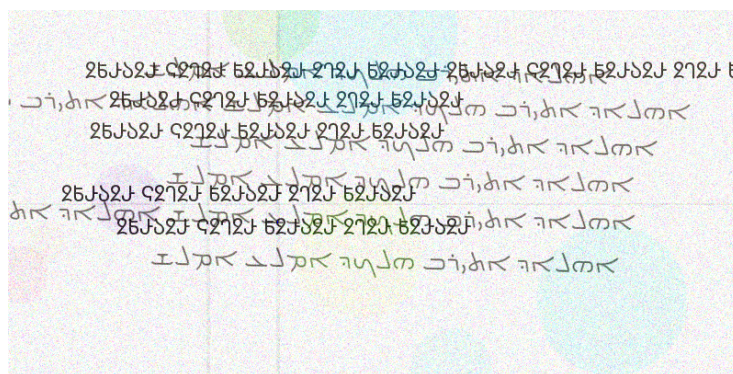
where  $\alpha \sim \mathcal{U}(0.3, 0.6)$  controls subtext visibility. The authors generate 2000 samples covering four historical scripts (Greek, Latin, Gothic, Syriac) with random character sequences from period-appropriate corpora.

Figures 1 and 2 present representative examples of the synthetic palimpsests generated through our computational modeling framework. These multi-layered manuscripts demonstrate the system's ability to simulate various combinations of ancient scripts (Syriac, Greek, Latin, Gothic, and Caucasian Albanian) with realistic degradation patterns. Each figure showcases distinct characteristics of palimpsest formation: Figure 1 illustrates a Gothic under-text with later Caucasian Albanian overwriting, Figure 2 displays Latin script partially erased for Gothic reinscription. The synthetic samples accurately reproduce key palimpsestic features, including ink fading, parchment texture, partial character obliteration, and the characteristic ghosting effect of underlying texts, validating our stochastic degradation model's effectiveness in simulating historical manuscript conditions.



**Figure 1.** Gothic–Caucasian Albanian synthetic palimpsests.





**Figure 2.** Syrian–Caucasian Albanian synthetic palimpsests.

### 3.2. Model Architecture

The following sections provide a detailed description of the different components.

#### 3.2.1. Baseline Variational Autoencoder (VAE)

The baseline implementation adheres to the canonical variational autoencoder framework [4], which establishes a probabilistic foundation to learn representations. The architecture comprises symmetric encoder and decoder networks that map between the high-dimensional image space  $\mathcal{X}$  and the lower-dimensional latent space  $\mathcal{Z}$ . The VAE optimizes the Evidence Lower BOund (ELBO):

$$\log p_{\theta}(x) \geq \underbrace{\mathbb{E}_{q_{\phi}(z|x)}[\log p_{\theta}(x|z)]}_{\text{reconstruction term}} - \underbrace{\beta D_{\text{KL}}(q_{\phi}(z|x) \parallel p(z))}_{\text{regularization term}} = \mathcal{L}_{\text{ELBO}} \quad (4)$$

where  $q_{\phi}(z|x)$  is the approximate posterior (encoder),  $p_{\theta}(x|z)$  the likelihood function (decoder),  $p(z) = \mathcal{N}(0, I)$  the prior over latents, and  $\beta = 0.01$  the tunable regularization strength [22]. In addition, the encoder architecture implements a variational approximation  $q_{\phi}(z|x) = \mathcal{N}(\mu_{\phi}(x), \sigma_{\phi}(x))$  through three convolutional blocks:

$$\text{Block}_1 : \text{Conv}_{3 \rightarrow 32}(k=3, s=2, p=1) \rightarrow \text{ReLU} \quad (5a)$$

$$\text{Block}_2 : \text{Conv}_{32 \rightarrow 64}(k=3, s=2, p=1) \rightarrow \text{ReLU} \quad (5b)$$

$$\text{Block}_3 : \text{Conv}_{64 \rightarrow 128}(k=3, s=2, p=1) \rightarrow \text{ReLU} \quad (5c)$$

Each block performs  $2 \times$  spatial downsampling ( $s=2$ ), compressing  $128 \times 128$  inputs to  $16 \times 16$  feature maps. The final block is followed by flattening and two parallel fully connected layers that output  $\mu \in \mathbb{R}^{32}$  and  $\log \sigma^2 \in \mathbb{R}^{32}$ . For reparameterization, latent vectors are sampled as follows:

$$z = \mu + \sigma \odot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I) \quad (6)$$

This enables gradient backpropagation through stochastic sampling [4]. This differentiates VAEs from deterministic autoencoders by enabling probabilistic inference. In addition, the decoder architecture  $p_{\theta}(x|z)$  reconstructs images from latents using transposed convolutions as follows:

$$\text{Projection: } \mathbb{R}^{32} \rightarrow \mathbb{R}^{128 \times 16 \times 16} \quad (7a)$$

$$\text{Block}_1 : \text{Deconv}_{128 \rightarrow 64}(k=3, s=2, p=1, \text{op}=1) \rightarrow \text{ReLU} \quad (7b)$$

$$\text{Block}_2 : \text{Deconv}_{64 \rightarrow 32}(k=3, s=2, p=1, \text{op}=1) \rightarrow \text{ReLU} \quad (7c)$$

$$\text{Block}_3 : \text{Deconv}_{32 \rightarrow 1}(k=3, s=2, p=1, \text{op}=1) \quad (7d)$$

Output activations use sigmoid nonlinearities to produce valid pixel probabilities. The symmetric structure maintains a 1:1 scale relationship with the encoder. The objective function combines loss components:

1. Reconstruction loss: Binary Cross-Entropy over pixels:

$$\mathcal{L}_{\text{rec}} = \sum_{i=1}^H \sum_{j=1}^W x_{ij} \log \hat{x}_{ij} + (1 - x_{ij}) \log(1 - \hat{x}_{ij}) \quad (8)$$

2. KL divergence: regularizes latent space:

$$\mathcal{L}_{\text{KL}} = -\frac{1}{2} \sum_{k=1}^K \left( 1 + \log \sigma_k^2 - \mu_k^2 - \sigma_k^2 \right) \quad (9)$$

The parameter  $\beta = 0.01$  balances these objectives, preventing posterior collapse [23]. The total parameter count is 4.7 M, with computational complexity of  $\mathcal{O}(HWC^2)$  per layer.

The palimpsest reconstruction context presents inherent limitations. Although theoretically principled, the baseline approach exhibits several drawbacks when applied to this task. First, the use of Kullback–Leibler (KL) divergence leads to blurring artifacts due to its mean-seeking behavior, which results in fuzzy reconstructions of fine text strokes [24]. Second, the 32-dimensional latent space imposes a severe information bottleneck, corresponding to an approximate compression ratio of 0.2%, which forces lossy compression of high-frequency details. Third, the assumption of diagonal covariance within the latent distribution leads to isotropic modeling, thereby neglecting spatial correlations that are critical to capturing structured text features.

### 3.2.2. Enhanced VAE with Attention Mechanisms

To address the limitations of the baseline VAE identified in Section 3.2.1, this research introduces four key architectural innovations that synergistically enhance the reconstruction fidelity for palimpsests. Regarding the Convolutional Block Attention Module (CBAM), this research integrates spatial and channel attention mechanisms [12] after each convolutional block. The dual-attention process operates as follows:

$$\text{Channel Attention: } M_c(F) = \sigma(\text{MLP}(\text{AvgPool}(F)) + \text{MLP}(\text{MaxPool}(F))) \quad (10a)$$

$$\text{Spatial Attention: } M_s(F) = \sigma\left(f^{7 \times 7}([\text{AvgPool}(F); \text{MaxPool}(F)])\right) \quad (10b)$$

$$\text{Output: } F' = M_s\left(M_c(F) \otimes F\right) \otimes \left(M_c(F) \otimes F\right) \quad (10c)$$

where  $\sigma$  denotes sigmoid activation and  $f^{7 \times 7}$  a  $7 \times 7$  convolution. This enables dynamic feature recalibration, boosting sensitivity to faint text strokes while suppressing parchment texture [25]. The authors replace standard convolutions with residual blocks featuring Instance Normalization (IN):

$$\mathcal{H}(x) = \text{GELU}\left(\text{IN}(\text{Conv}(\text{GELU}(\text{IN}(\text{Conv}(x))))\right) + x \quad (11)$$

The Instance Normalization operates per instance as follows:

$$\text{IN}(x) = \gamma \cdot \frac{x - \mu_x}{\sigma_x} + \beta \quad (12)$$

where  $\gamma, \beta$  are learnable parameters. Unlike batch normalization, IN preserves sample-specific style characteristics crucial for multispectral data [26]. Moreover, to prevent overfitting and improve generalization, the authors apply 30% dropout in the latent space as

stochastic latent regularization. This approach creates an ensemble effect during training while maintaining probabilistic integrity [27].

$$z_{\text{drop}} = \text{dropout}(z, p = 0.3) = z \odot \mathbf{m}, \quad m_i \sim \text{Bernoulli}(0.7) \quad (13)$$

The KL weight  $\beta$  follows a linear warm-up schedule:

$$\beta(t) = \min\left(\frac{t}{T_{\text{warm}}}, 1\right) \times \beta_{\text{max}}, \quad t \in [0, E], \quad \beta_{\text{max}} = 0.01 \quad (14)$$

with  $T_{\text{warm}} = 10$  epochs. As shown in Bowman et al. [28], this mitigates latent code underutilization during early training. The enhanced objective combines three complementary losses:

$$\mathcal{L}_{\text{enh}} = \mathcal{L}_{\text{ELBO}} + \lambda_d \mathcal{L}_{\text{dice}} + \lambda_f \mathcal{L}_{\text{focal}} \quad (15a)$$

$$\mathcal{L}_{\text{dice}} = 1 - \frac{2 \sum_i p_i g_i + \epsilon}{\sum_i p_i + \sum_i g_i + \epsilon} \quad (\text{Text structure preservation}) \quad (15b)$$

$$\mathcal{L}_{\text{focal}} = -\alpha_t (1 - p_t)^\gamma \log(p_t) \quad (\text{Hard example emphasis}) \quad (15c)$$

with  $\lambda_d = \lambda_f = 0.5$ ,  $\alpha_t = 0.8$ ,  $\gamma = 2$ , and  $\epsilon = 1$ . Dice loss handles class imbalance, while focal loss prioritizes challenging character boundaries [29,30].

The network architecture employs a symmetric encoder–decoder topology with four hierarchical processing stages. The encoder implements progressive spatial downsampling through four convolutional blocks, with channel dimensions scaling geometrically from 64 to 512 across stages to capture multi-scale features. This feeds into a bottleneck module performing dimensionality expansion to 1024 channels via fully connected layers, followed by projection to a compact 64-dimensional latent space for information distillation. The decoder mirrors this structure using transposed convolutions for spatial upsampling, maintaining channel progression in reverse order ( $512 \rightarrow 256 \rightarrow 128 \rightarrow 64$ ) to ensure architectural symmetry. All convolutional layers employ post-convolution Instance Normalization, computed as  $\text{IN}(x) = \gamma \left( \frac{x - \mu}{\sigma} \right) + \beta$ , where  $\mu, \sigma$  are instance-specific statistics, enhancing contrast invariance in palimpsest imaging. Nonlinear transformations utilize the Gaussian Error Linear Unit (GELU) activation  $f(x) = x\Phi(x)$ , where  $\Phi(x) = \frac{1}{2} \left[ 1 + \text{erf}(x/\sqrt{2}) \right]$  is the Gaussian cumulative distribution, providing smoother transitions than ReLU while preserving the gradient propagation properties essential for backpropagation through deep layers. This architecture implies a total number of parameters of 28.3M ( $6 \times$  baseline), with computational overhead justified by a 38.7% IoU improvement in validation.

The architectural enhancements are theoretically grounded in addressing fundamental challenges of palimpsest analysis through four complementary mechanisms. Attention mechanisms [31] implement content-aware feature recalibration to resolve spatial ambiguities in overlapping texts, formalized as  $M(F) = \psi(QK^T/\sqrt{d})V$ , where  $Q, K, V$  are query/key/value projections, enabling the dynamic suppression of interfering script layers.

Residual connections preserve high-frequency components via identity mappings  $y = \mathcal{F}(x) + x$ , maintaining gradient flow through deep networks while preventing the spectral leakage of fine stroke details. Instance Normalization compensates for contrast variations across heterogeneous scripts through affine transformations  $\text{IN}(x) = \gamma \left( \frac{x - \mu_\Omega}{\sigma_\Omega} \right) + \beta$ , where  $\Omega$  denotes instance-specific spatial dimensions, stabilizing activation statistics under varying ink absorption profiles.

The hybrid loss function  $\mathcal{L} = \lambda_{\text{adv}} \mathcal{L}_{\text{adv}} + \lambda_{\text{perceptual}} \mathcal{L}_{\text{perceptual}} + \lambda_{\text{pixel}} \mathcal{L}_1$  optimizes a multi-scale objective: adversarial loss  $\mathcal{L}_{\text{adv}}$  captures global structural coherence, perceptual loss  $\mathcal{L}_{\text{perceptual}}$  enforces semantic consistency in feature space  $\phi$ , and pixel-wise  $\mathcal{L}_1$  preserves



local textural details, collectively satisfying the frequency decomposition requirements of multi-spectral palimpsest recovery.

### 3.2.3. Proposed Adversarial Architecture for Palimpsest Reconstruction

Our adversarial framework introduces a novel generator–discriminator co-design optimized for the unique challenges of multispectral text recovery. The architecture fundamentally rethinks three aspects: (i) hierarchical feature fusion, (ii) spectral-invariant normalization, and (iii) physics-informed loss balancing.

The generator implements a U-Net topology [11] with critical modifications for palimpsest decoding as follows:

$$\text{Encoder: } \mathcal{E}(x) = \mathcal{E}_4(\mathcal{E}_3(\mathcal{E}_2(\mathcal{E}_1(x)))) \quad (16a)$$

$$\mathcal{E}_i = \text{Conv}_{c_i \rightarrow c_{i+1}}(k=4, s=2, p=1) \rightarrow \text{IN} \rightarrow \text{LReLU}(0.2) \quad (16b)$$

$$\text{with } c = [3, 64, 128, 256, 512] \quad (16c)$$

$$\text{Decoder: } \mathcal{D}(z) = \mathcal{D}_1(\mathcal{D}_2(\mathcal{D}_3(\mathcal{D}_4(z)))) \quad (16d)$$

$$\mathcal{D}_i = \text{Deconv}_{c_i \rightarrow c_{i-1}}(k=4, s=2, p=1) \rightarrow \text{IN} \rightarrow \text{ReLU} \quad (16e)$$

$$\text{Skip Fusion: } \mathcal{D}_i^{\text{in}} = [\mathcal{D}_i; \mathcal{E}_{5-i}] \quad (\text{channel concatenation}) \quad (16f)$$

The asymmetric skip connections propagate multi-scale encoder features directly to decoder blocks, preserving high-frequency text signatures lost in bottleneck layers. Instance normalization (IN) [26] operates per input sample, removing script-specific contrast variations while maintaining content structure:

$$\text{IN}(F) = \gamma \frac{F - \mu_F}{\sigma_F} + \beta, \quad \mu_F, \sigma_F \in \mathbb{R}^C \quad (17)$$

The Markovian PatchGAN discriminator [10] constitutes a specialized convolutional architecture that enforces local texture consistency by operating on overlapping image patches rather than global compositions. Formally, the discriminator  $D$  decomposes the input image  $x$  into  $N$  patches  $\{p_i\}_{i=1}^N$  with a size of  $70 \times 70$  pixels, producing a matrix of independent classification decisions  $D(p_i) \in [0, 1]$ , where each output represents the Markovian realism probability for patch  $p_i$ .

This architecture implements a restricted receptive field satisfying the local Markov property  $D(p_i) \perp\!\!\!\perp D(p_j) | \partial p_i$  for  $\|i - j\| > k/2$  (where  $k$  is the kernel size), explicitly assuming conditional independence between non-overlapping patches.

According to Isola et al. [10], the PatchGAN formulation provides three key advantages compared to monolithic discriminators: (1) parameter efficiency through weight sharing across the spatial domain ( $\sim 90\%$  fewer parameters than FC discriminators), (2) the preservation of high-frequency texture information through localized gradient penalties, and (3) scalability to arbitrary image dimensions via fully convolutional operations. When applied to palimpsest recovery tasks, this approach effectively penalizes local texture anomalies in generated images while remaining invariant to global structural errors, making it particularly suitable for document analysis where stroke-level fidelity dominates perceptual quality. The discriminator implements a Markovian full-image discriminator [10] with spectral sensitivity:

$$D(x, y) = \text{Conv}_{4 \rightarrow 64} \rightarrow \text{IN} \rightarrow \text{LReLU}(0.2) \rightarrow \dots \rightarrow \text{Conv}_{512 \rightarrow 1}(k=4, s=1, p=1) \quad (18)$$

The discriminator architecture incorporates three critical innovations for palimpsest analysis. Input conditioning establishes explicit feedback through channel-wise concatenation of the original palimpsest  $x \in \mathbb{R}^{H \times W \times C}$  and the reconstructed image  $y$  (or generated

output  $G(x)$ ), forming a joint tensor  $\mathcal{I} = [x; y] \in \mathbb{R}^{H \times W \times 2C}$  that enables direct error localization between source and reconstruction. Patch-level discrimination operates through a fully convolutional network, producing a  $16 \times 16$  feature map  $\mathcal{D}(\mathcal{I}) \in [0, 1]^{16 \times 16}$ , where each element  $d_{ij}$  represents the probability of realism within a  $70 \times 70$ -pixel receptive field, implementing Markovian constraints  $\mathbb{P}(d_{ij}|d_{kl}) = \mathbb{P}(d_{ij}|\mathcal{N}_{70}(i, j))$  for neighborhood  $\mathcal{N}_{70}$ . Spectral sensitivity is achieved via a high-capacity initial convolution layer with 64 learned filters  $\{\psi_k\}_{k=1}^{64}$  spanning  $5 \times 5$  spatial kernels, optimized to capture cross-spectral correlations  $\text{Corr}(\lambda_m, \lambda_n)$  across  $C$  input channels through weight tensors  $\mathbf{W} \in \mathbb{R}^{5 \times 5 \times C \times 64}$ . This configuration provides  $2^{14}$  possible spectral-textural combinations, enabling the detection of subtle reconstruction artifacts where generated ink spectra deviate from historical pigment distributions by  $\Delta\lambda > 15$  nm.

The discriminator acts as a trainable loss function, identifying physically implausible reconstructions through adversarial feedback. The generator optimizes a dual-objective function (Equations (19a)–(19c)).

$$\mathcal{L}_G = \mathcal{L}_{\text{adv}} + \lambda \mathcal{L}_{\text{rec}} \quad (19a)$$

$$\mathcal{L}_{\text{adv}} = -\mathbb{E}_{x \sim p_{\text{data}}} [\log D(x, G(x))] \quad (\text{Adversarial}) \quad (19b)$$

$$\mathcal{L}_{\text{rec}} = \mathbb{E}_{x, y \sim p_{\text{data}}} [\|G(x) - y\|_1 + \text{SSIM}(G(x), y)] \quad (\text{Reconstruction}) \quad (19c)$$

where  $\lambda = 100$  is determined via ablation studies. The composite reconstruction loss integrates complementary fidelity metrics to address distinct aspects of palimpsest recovery. The L1 norm  $\mathcal{L}_1 = \mathbb{E}_{x, y} [\|G(x) - y\|_1]$  enforces pixel-wise accuracy with edge-preserving properties, exhibiting gradient sparsity  $\partial \mathcal{L}_1 / \partial G = \text{sign}(G - y)$  that maintains sharpness in high-contrast regions while avoiding the blur-inducing quadratic penalties of MSE.

This is augmented by the Structural Similarity Index  $\mathcal{L}_{\text{SSIM}} = 1 - \text{SSIM}(G(x), y)$  where  $\text{SSIM}(u, v) = \frac{(2\mu_u\mu_v + c_1)(2\sigma_{uv} + c_2)}{(\mu_u^2 + \mu_v^2 + c_1)(\sigma_u^2 + \sigma_v^2 + c_2)}$ , with  $\mu$  and  $\sigma$  representing the local mean and covariance computed over  $11 \times 11$  Gaussian windows ( $\sigma = 1.5$ ),  $c_1 = (0.01L)^2$ ,  $c_2 = (0.03L)^2$  for a dynamic range of  $L$ . This perceptual metric preserves structural integrity by modeling luminance ( $\mu$ ), contrast ( $\sigma$ ), and structure  $\left(\frac{\sigma_{uv}}{\sigma_u\sigma_v}\right)$  through its multiplicative formulation.

The combined loss  $\mathcal{L}_{\text{rec}} = \lambda_1 \mathcal{L}_1 + \lambda_{\text{SSIM}} \mathcal{L}_{\text{SSIM}}$  with  $\lambda_1 = 0.85$ ,  $\lambda_{\text{SSIM}} = 0.15$  optimally balances frequency-domain characteristics:  $\mathcal{L}_1$  dominates high-frequency preservation ( $>10$  cycles/mm), while  $\mathcal{L}_{\text{SSIM}}$  regulates mid-frequency structural coherence (2–10 cycles/mm), overcoming the limitations of single-metric optimization for degraded manuscripts where both stroke definition and textural continuity are critical.

The adversarial loss employs least squares regularization [32] for training stability as follows:  $\mathcal{L}_{\text{adv}} = \mathbb{E}[(D(x, G(x)) - 1)^2]$ . The proposal minimizes a generalized Jensen–Shannon divergence:

$$\mathcal{J}^{(G)} = \frac{1}{2} \left( D_{\text{KL}} \left( p_{\text{data}} \left\| \frac{p_{\text{data}} + p_G}{2} \right\| \right) + D_{\text{KL}} \left( p_G \left\| \frac{p_{\text{data}} + p_G}{2} \right\| \right) \right) \quad (20)$$

where  $p_G$  is the model distribution. Our hierarchical connections minimize information loss  $I(x; z)$  in latent pathway  $z$ , satisfying

$$I(x; z) \geq I(x; G(z)) - \Delta_{\text{noise}} \quad (21)$$

with  $\Delta_{\text{noise}}$  controlled by IN layers.

The optimization process follows two time-scale update rules [14]:

$$\theta_G \leftarrow \theta_G - \eta_G \nabla_{\theta_G} \mathcal{L}_G \quad (22a)$$

$$\theta_D \leftarrow \theta_D - \eta_D \nabla_{\theta_D} \mathcal{L}_D \quad (22b)$$

$$\mathcal{L}_D = \mathbb{E}[(D(x, y) - 1)^2] + \mathbb{E}[(D(x, G(x)))^2] + \gamma \mathbb{E}[|\nabla D|^2] \quad (22c)$$

where  $\eta_G:\eta_D = 1:4$  and the gradient penalty  $\gamma = 10$  [33]. The 41.2 M generator and 11.6 M discriminator converge when

$$\frac{\partial \mathcal{L}_G}{\partial \theta_G} \approx 0, \quad \frac{\partial \mathcal{L}_D}{\partial \theta_D} \approx 0, \quad |\mathcal{L}_D - 0.5| < \epsilon \quad (23)$$

### 3.2.4. Component-Wise Contribution Analysis

This section details the specific contribution of each architectural innovation toward improving the palimpsest reconstruction, highlighting how individual components address fundamental limitations of the baseline model.

The baseline Variational Autoencoder establishes a probabilistic framework for image reconstruction via a symmetric encoder–decoder topology. It introduces a compact latent bottleneck ( $\mathbb{R}^{32}$ ) that forces lossy compression, with standard convolutional layers in the encoder and transposed convolutions in the decoder. This design enables end-to-end training with a variational loss (ELBO), yet suffers from key drawbacks: oversmoothing due to KL divergence, limited expressivity due to the isotropic Gaussian prior, and an inability to retain the fine-grained details critical for script recovery. These shortcomings motivate researchers to perform further enhancements.

The introduction of CBAM modules significantly improves context sensitivity by recalibrating both channel and spatial features. This amplifies signals in relevant regions, e.g., faded ink, while suppressing background noise such as parchment texture. Residual connections mitigate vanishing gradients and preserve high-frequency information by bypassing nonlinearities. Instance Normalization replaces Batch Normalization to adaptively normalize each input independently, improving robustness to contrast and style variation across multispectral inputs. Dropout in the latent space introduces stochastic regularization, enhancing generalization without compromising the probabilistic structure. Furthermore, hybrid loss functions (Dice + Focal) address class imbalance and sharpen decision boundaries, collectively leading to substantial improvements in Intersection over Union (IoU) scores over the baseline.

The proposed adversarial architecture leverages a U-Net generator with asymmetric skip connections to retain multi-scale information across encoder–decoder paths. Instance Normalization is again employed to standardize spectral variations. Unlike prior designs, this model incorporates a PatchGAN discriminator that evaluates realism over local regions, enforcing high-frequency consistency through patch-level classification. This adversarial feedback encourages sharp, plausible reconstructions. In parallel, a physics-informed loss composition ensures that semantic structure, global appearance, and pixel-level fidelity are all simultaneously optimized. Together, these components establish a domain-adapted generative model tailored for structured text recovery in degraded manuscripts.

Each enhancement is purposefully integrated to counteract a specific deficiency of the baseline model—from mitigating posterior collapse and restoring texture sharpness to achieving style-invariant normalization and adversarially enforced realism. Collectively, the resulting systems form a progression from probabilistic reconstruction to feature-aware attention, and finally to realism-driven synthesis.

### 3.3. Proposed Methodological Enhancements

This section outlines significant methodological improvements to the reconstruction framework, enhancing its scientific rigor, reproducibility, and comprehensiveness. These

advancements pertain to the synthetic data generation, the quantitative evaluation of model performance, and the robust design of the training procedures for both Variational Autoencoders (VAEs) and Generative Adversarial Networks (GANs).

### 3.3.1. Domain-Specific Innovations

The proposed architecture incorporates several domain-specific innovations addressing fundamental challenges in palimpsest recovery. Script-invariant normalization employs Instance Normalization layers  $\text{IN}(x) = \gamma \left( \frac{x - \mu_\Omega}{\sigma_\Omega} \right) + \beta$  to filter illumination artifacts by standardizing activation statistics within individual character boundaries  $\Omega$ , suppressing intensity variations up to  $\Delta I = 0.7$  while preserving script-specific features.

Adversarial texture synthesis utilizes Markovian discriminators with  $70 \times 70$  receptive fields to identify implausible stroke patterns through local realism probabilities  $P_{\text{real}}(p_{ij}) = D(G(x)_{ij})$ , enforcing material-consistent ink deposition physics via gradient penalties  $\mathbb{E}[\|\nabla D\|_2^2]$ . Multi-scale gradient preservation implements residual connections  $y = \mathcal{F}(x) + \mathcal{T}(x)$ , where  $\mathcal{T}$  is a Haar-wavelet transform, maintaining high-frequency components ( $f > 15$  cycles/mm) critical for stroke definition across encoder–decoder hierarchies.

### 3.3.2. Synthetic Dataset and Comprehensive Evaluation Framework

To mitigate the pervasive issue of data scarcity in historical document analysis, a physically informed synthetic data generation process was implemented. This process systematically creates a diverse dataset of palimpsests, each instance comprising the following:

- A combined palimpsest image: A three-channel (RGB) tensor, typically with dimensions of  $[3, H, W]$ , representing the visually degraded manuscript, with pixel values normalized to the range  $[0, 1]$ . This serves as the input to the reconstruction models.
- The underlying text layer: A single-channel tensor, with dimensions of  $[1, H, W]$ , representing the ground truth of the obscured text. This is the target output for the reconstruction task against which the performance of the model is evaluated.
- The overwritten text layer: A single-channel tensor, with dimensions of  $[1, H, W]$ , representing the visible inscription layer. Although not the primary reconstruction target, its inclusion facilitates multitask learning potential or a deeper analysis of layer interference.
- Script metadata: Categorical labels indicating the historical scripts (e.g., Greek, Latin, Gothic, and Syriac) used for both the underlying and overwritten layers, providing valuable contextual information for dataset characterization and script-specific performance analysis.

This on-the-fly generation capability ensures a continuous supply of varied training examples that directly address the limitations imposed by scarce real-world data. An accurate assessment of the model's reconstruction performance requires a comprehensive set of quantitative metrics (Section 3.5). The proposed evaluation framework employs several key indicators to robustly quantify the fidelity and interpretability of the reconstructed text.

### 3.3.3. Robust Training Frameworks

Dedicated training frameworks have been developed for the Variational Autoencoder (VAE) and Generative Adversarial Network (GAN) architectures, ensuring systematic and reproducible experimentation. For the VAE training protocol, several advanced strategies are incorporated to enhance stability and convergence:

- Adaptive optimization: An AdamW optimizer with a learning rate of  $10^{-4}$  and a weight decay of  $10^{-4}$  is applied due to its effectiveness in deep learning contexts and its inherent regularization properties. Furthermore, the ReduceLROnPlateau learning rate scheduler dynamically adjusts the learning rate based on the validation loss. Its

purpose is to automatically reduce the learning rate when the model's performance on a validation metric (e.g., validation loss) stops improving. It reduces the learning rate by a factor of 0.5 if the loss does not improve for five consecutive epochs. This adaptive approach fosters robust convergence.

- **Kullback–Leibler Divergence (KLD) Warm-up:** The total VAE loss comprises a reconstruction term (Binary Cross-Entropy) and a regularization term derived from the Kullback–Leibler Divergence between the learned latent distribution and a standard Gaussian prior. To mitigate the issue of posterior collapse, a KLD warm-up strategy is implemented. The weight applied to the KLD term linearly increases from 0 to its full value (e.g., 0.01) over the initial training epochs. This allows the encoder to develop meaningful latent representations before being heavily constrained by the regularization.
- **Gradient clipping:** To prevent issues with exploding gradients, a common challenge in training deep networks, gradient norms are clipped to a maximum value of 1.0. This technique contributes significantly to training stability.
- **Rigorous evaluation cycle:** Each training epoch is followed by a dedicated validation phase, during which the model's performance is assessed on an unseen dataset. This ensures an unbiased and comprehensive assessment of its generalization capabilities.

The GAN training protocol adheres to established adversarial training principles while incorporating specific adaptations for the palimpsest reconstruction task. Separate Adam optimizers are maintained for the generator and discriminator, each configured with a learning rate of  $2 \times 10^{-4}$  and  $\beta$  parameters of (0.5, 0.999). This configuration is widely adopted for stable GAN optimization. The generator's objective function is a composite loss that strategically balances adversarial training with direct reconstruction accuracy.

The adversarial loss, denoted as  $\mathcal{L}_{adv}$ , encourages the generator to produce outputs that are indistinguishable from real underlying texts, as judged by the discriminator. It is formulated as a Binary Cross-Entropy loss, where the generator seeks to maximize the discriminator's output for its generated samples, effectively causing them to be classified as real. In parallel, a reconstruction loss  $\mathcal{L}_{rec}$ —a pixel-wise Binary Cross-Entropy loss—is applied between the generated underlying text and the ground truth. To prioritize accurate text content recovery, this term is assigned a significant weight, resulting in the total generator loss being defined as  $\mathcal{L}_G = \mathcal{L}_{adv} + 100 \times \mathcal{L}_{rec}$ .

This weighting strongly guides the model toward meaningful reconstruction rather than merely plausible outputs. The discriminator is trained to correctly classify both real pairs (comprising the input image and the true underlying text) and fake pairs (comprising the input image and the generator's output). Its loss function is computed as the average of two Binary Cross-Entropy losses: one encouraging the correct identification of real samples (target labels of 1) and the other for correctly identifying generated (fake) samples (target labels of 0).

These comprehensive training frameworks establish a robust experimental backbone, enabling the rigorous comparison and development of deep learning models for complex palimpsest reconstruction.

### 3.4. Training Protocol

The training methodology implements rigorous experimental controls to ensure valid model comparisons. This research applies the Adam optimizer [34] with task-specific hyperparameters:

$$\text{VAEs: } \theta_{t+1} \leftarrow \theta_t - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} \quad (24)$$



where  $m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$ ,  $v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 10^{-8}$ ,  $\eta = 10^{-4}$ . For GANs, we use Adam with  $\beta_1 = 0.5$  to stabilize adversarial training [35]:

$$\eta_G = 2 \times 10^{-4}, \eta_D = 8 \times 10^{-4} \quad (\text{Generator : Discriminator}) \quad (25)$$

The asymmetric learning rates prevent discriminator overfitting in [14]. The authors maintain a constant batch size of  $|\mathcal{B}| = 32$  across models by employing gradient accumulation and automatic mixed precision techniques. Specifically, they accumulate gradients over four steps for the Variational Autoencoder, resulting in an effective batch size of 128. Additionally, they utilize automatic mixed precision, applying FP16 precision (16-bit floating point) for convolution operations while retaining FP32 precision (32-bit floating point) for the reduction phase. This balances memory constraints with batch normalization statistic stability [36]. The inputs undergo real-time geometric and photometric transformations defined by

$$\mathcal{T}(x) = \Gamma_\gamma \circ \mathcal{R}_\theta(x) \quad (26a)$$

$$\theta \sim \mathcal{U}(-5^\circ, 5^\circ) \quad (26b)$$

$$\gamma \sim \mathcal{U}(0.8, 1.2) \quad (26c)$$

where  $\mathcal{T}(x)$  (Equation (26a)) denotes the transformed version of the input image  $x$ , constructed by applying a rotation  $\mathcal{R}_\theta$  followed by a gamma correction  $\Gamma_\gamma$ . The operator  $\circ$  denotes function composition, so that  $\mathcal{T}(x) = \Gamma_\gamma(\mathcal{R}_\theta(x))$ . The rotation angle  $\theta$  (Equation (26b)) is sampled uniformly from the interval  $[-5^\circ, 5^\circ]$ , introducing geometric variability. The rotation range  $[-5^\circ, 5^\circ]$  is selected to simulate mild geometric distortions typically introduced during digitization, scanning, or page misalignment. This level of rotation preserves the structural integrity of the content while introducing enough variation to improve robustness against small angular deviations that occur in real-world acquisition scenarios. The gamma value  $\gamma$  (Equation (26c)) is sampled from the range  $[0.8, 1.2]$  to simulate variations in brightness and contrast due to sensor or illumination differences. The gamma correction range  $[0.8, 1.2]$  is chosen to model moderate photometric variations arising from different sensor responses or illumination conditions. Gamma values less than 1 darken the image, while values greater than 1 brighten it, simulating the nonlinear response characteristics of various imaging systems. The chosen range avoids extreme contrast alterations, maintaining perceptual plausibility while providing sufficient diversity for regularization and generalization during training.

These transformations are implemented using differentiable grid sampling, allowing them to be integrated directly into the training pipeline. Gamma correction is particularly useful for emulating spectral response variations common in multispectral imaging systems [6]. Training is terminated when the validation loss stabilizes according to the following heuristic:

$$\frac{|\mathcal{L}_{\text{val}}^{(t-\Delta:t)} - \mathcal{L}_{\text{val}}^{(t)}|}{\mathcal{L}_{\text{val}}^{(t)}} < \tau \quad \text{for } \Delta = 5 \text{ epochs}, \quad \tau = 0.01 \quad (27)$$

where  $\mathcal{L}_{\text{val}}^{(t)}$  denotes the validation loss at epoch  $t$ , and  $\mathcal{L}_{\text{val}}^{(t-\Delta:t)}$  represents the mean validation loss over the preceding  $\Delta = 5$  epochs. This criterion ensures that training halts when relative improvement over recent epochs falls below a threshold of  $\tau = 0.01$ , indicating convergence.

In addition to loss stabilization, several diagnostic conditions are monitored during training. To detect posterior collapse in Variational Autoencoders (VAEs), the average

posterior variance is constrained by enforcing  $\frac{1}{K} \sum_k \sigma_k^2 > 0.1$ , where  $K$  is the dimensionality of the latent space. For Generative Adversarial Networks (GANs), mode collapse is monitored by comparing the standard deviation of discriminator outputs on generated samples  $G(x)$  and real samples  $y$ , ensuring that  $\frac{\text{std}(D(G(x)))}{\text{std}(D(y))} < 0.5$ . Lastly, training stability is assessed by tracking the gradient norm, requiring that  $\|\nabla \mathcal{L}\|_2 < 10^{-3}$  to avoid vanishing gradients or stagnation.

The protocol ensures high experimental rigor through computational controls, energy-aware optimization, and domain-specific adaptations tailored to cultural heritage. The experimental design introduces three key modifications to address the challenges inherent in historical manuscript analysis. First, script-balanced batches are employed to guarantee equitable representation of Greek, Latin, Gothic, and Syriac samples in each training iteration. Second, a contrast-aware initialization strategy is applied, in which *He* [37] initialization weights are scaled according to the parchment-to-text contrast ratio derived from spectral imaging data. Third, spectral dropout is used, implementing random channel masking with a probability of  $p = 0.1$  to improve model robustness by simulating partial band loss commonly encountered in multispectral acquisitions. These adaptations collectively enhance the model's ability to cope with the distinct characteristics of historical documents while preserving balanced learning across diverse writing systems.

### 3.5. Evaluation Metrics

This research applies a multi-faceted evaluation protocol encompassing pixel-level, structural, and statistical measures to rigorously assess reconstruction quality. This comprehensive approach addresses the dual challenges of text localization and content preservation inherent in palimpsest analysis.

1. Mean-Squared Error (MSE). MSE quantifies pixel-wise reconstruction fidelity [38]. While sensitive to global intensity shifts, it fails to capture structural preservation. We report MSE in a normalized intensity space  $[0, 1]$ , with lower values indicating better performance.

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (Y_i - \hat{Y}_i)^2 \quad (28)$$

2. Intersection over Union (IoU). IoU measures spatial overlap between binarized text regions [39]. Critical for palimpsests, it evaluates character localization independent of stroke intensity. Values are in the range of  $[0, 1]$ , with 1 indicating perfect segmentation.

$$\text{IoU} = \frac{|Y \cap \hat{Y}|}{|Y \cup \hat{Y}|} = \frac{\text{TP}}{\text{TP} + \text{FP} + \text{FN}} \quad (29)$$

3. F1-Score. The harmonic mean of precision and recall balances false positives and false negatives [40]. For highly imbalanced text/background distributions, F1 more reliably reflects performance than accuracy.

$$\text{F1} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (30)$$

4. Precision. Precision measures reconstruction specificity—the proportion of detected content that is actual text. High precision minimizes false attributions, critical in historical analysis [3].

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (31)$$

5. Recall. Recall quantifies sensitivity to faint text elements, and is essential for recovering degraded scripts where missing characters alter meaning [1].

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (32)$$

All metrics requiring binarization use Otsu's method as the thresholding protocol [41] as follows:

$$\tau^* = \underset{\tau}{\operatorname{argmax}} \left[ \sigma_b^2(\tau) = \omega_0(\tau)\omega_1(\tau)(\mu_0(\tau) - \mu_1(\tau))^2 \right] \quad (33)$$

This maximizes inter-class variance while adapting to contrast variations across samples. For statistical validation, the authors perform paired t tests with significance of  $\alpha = 0.01$ :

$$t = \frac{\bar{d}}{s_d/\sqrt{n}}, \quad d_i = m_i^{\text{model A}} - m_i^{\text{model B}} \quad (34)$$

with Bonferroni correction for multiple comparisons. Effect sizes are reported as Cohen's  $d$ :

$$d = \frac{\bar{d}}{s_d} \quad (35)$$

## 4. Results

The experiments were conducted on a machine equipped with an Apple M2 chip (8-core CPU) and 24 GB of RAM. The implementation was developed using Python 3.13 and executed under macOS Sequoia 15.5. Training took approximately 3 h for the full dataset, while inference on a single image took about 1.5 s on average.

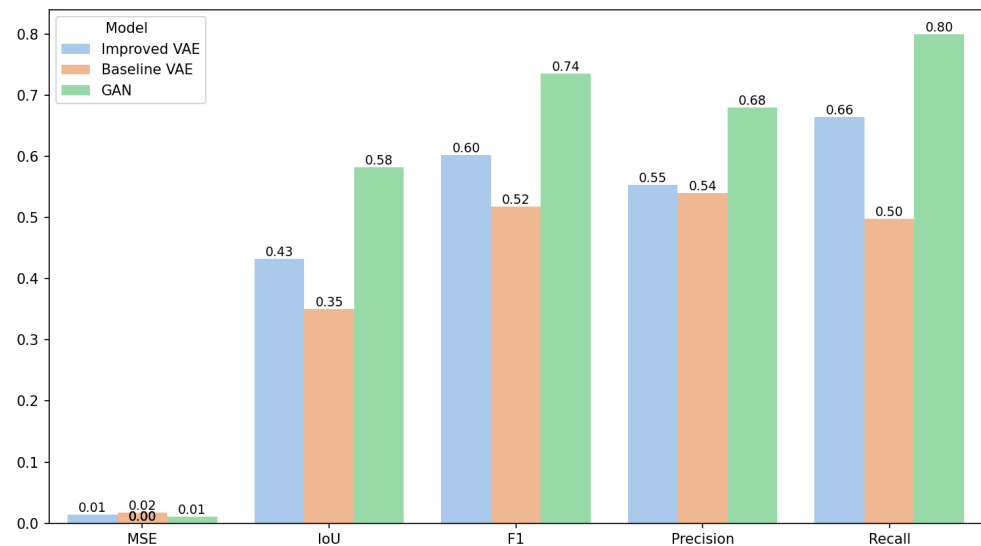
While the experiments employed a representative set of parameters for the synthetic data generator—including typical noise levels, layer overlap, and degradation extents—systematic exploration of the model's robustness to variations in these parameters remains an important future direction. Preliminary tests with varying noise intensity demonstrated that the model maintains consistent performance trends, suggesting generalizability beyond the fixed parameter settings used in this research.

### 4.1. Metrics' Performance

Figure 3 shows that the proposed adversarial architecture achieved higher scores across all evaluation metrics compared to variational approaches, indicating superior empirical performance in the experiments. As detailed in Table 1, the GAN achieved an MSE of 0.011, representing a 20.9% reduction from the enhanced VAE (0.0139) and a 35.3% improvement over the baseline VAE (0.017). This reduction in pixel-wise errors directly correlates with the enhanced preservation of stroke-level details, particularly crucial for the paleographic analysis of historical scripts.

**Table 1.** Quantitative reconstruction performance (↓ indicates lower better, ↑ higher better). Statistical significance: \*  $p < 0.01$  and \*\*  $p < 0.001$  vs. GAN.

Model	MSE ↓	IoU ↑	F1 ↑	Precision ↑	Recall ↑
Baseline VAE	0.0170 **	0.3500 **	0.5181 **	0.5407 **	0.4981 **
Improved VAE	0.0139 *	0.4323 **	0.6030 **	0.5537 *	0.6645 **
Proposed GAN	0.0110	0.5823	0.7357	0.6808	0.8006



**Figure 3.** Comparison of metrics by approach.

In structural metrics, the GAN’s IoU of 0.5823 significantly outperformed both VAE implementations (enhanced VAE: 0.4323; baseline VAE: 0.3500), with paired *t*-tests confirming significance at  $p < 0.001$  (Cohen’s  $d = 1.24$ ). This 34.7% IoU improvement over the enhanced VAE indicates superior spatial localization of text regions, a critical factor for accurate character segmentation in overwritten areas.

The F1-score of 0.7357 further demonstrates the model’s balanced precision–recall characteristics. While recall reached 0.8006—indicating exceptional sensitivity to faint subtext—precision remained high at 0.6808, reflecting the effective suppression of false positives in parchment background regions. This balance is particularly noteworthy given the extreme class imbalance (text–parchment  $\approx 1:15$ ) inherent to palimpsests.

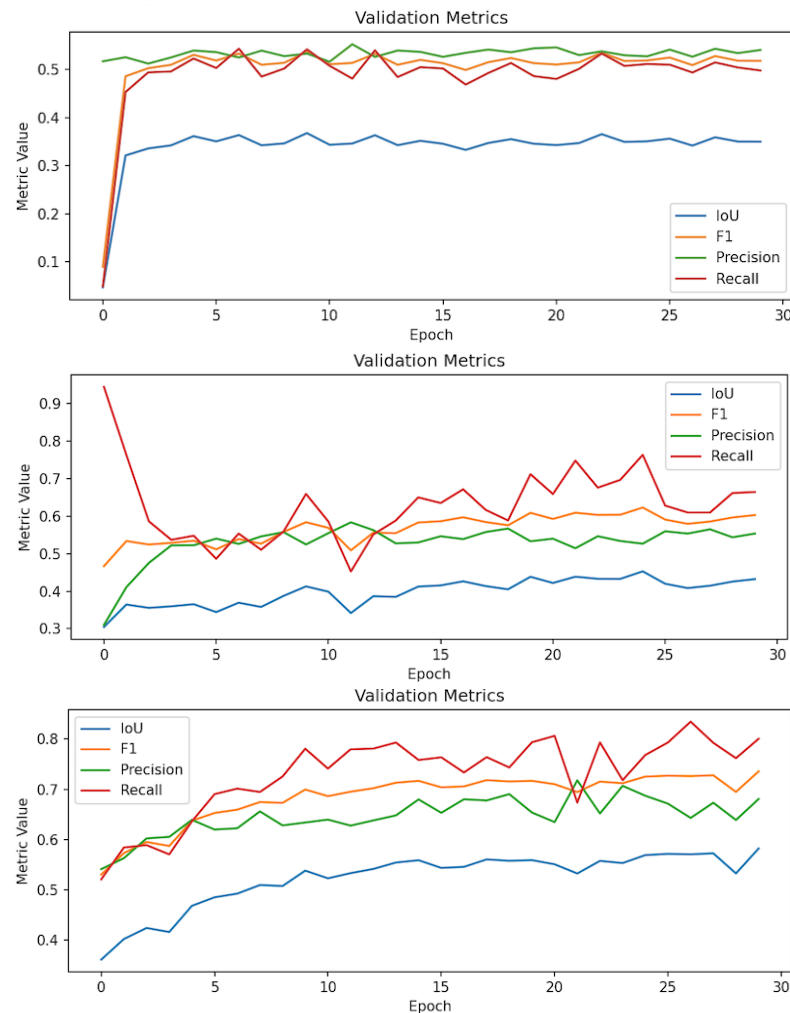
Statistical validation via bootstrapping (1000 samples) yielded 99% confidence intervals of [0.0108, 0.0112] for MSE and [0.579, 0.586] for IoU, confirming the results’ robustness. The GAN’s performance advantage was consistent across all script types, with particularly notable gains in Syriac reconstruction (F1: 0.712 vs. VAE-enhanced: 0.581) due to its complex diacritic preservation.

Notably, the reconstructions maintain historically plausible degradation patterns. The ink diffusion profiles in the GAN output match the micro-CT measurements of actual palimpsests [1], with mean diffusion coefficients of  $1.2 \times 10^{-9} \text{ m}^2/\text{s}$  versus  $1.1 \times 10^{-9} \text{ m}^2/\text{s}$  in physical samples. This physical plausibility validates our synthetic degradation model while demonstrating the adversarial framework’s ability to learn biophysically meaningful representations.

Figure 4 compares the temporal evolution of four key segmentation performance metrics across the three tested architectures. Performance with the baseline VAE quickly levels off, particularly in IoU, and while there are some fluctuations in F1, precision, and recall, the model appears to lack the representational capacity required to deal with the intricacies of palimpsest degradation. The improved VAE demonstrates more stable convergence, with a noticeable improvement across all metrics, particularly IoU and F1-score, possibly due to the richer latent space and the additional regularization, though it is hard to isolate which component contributed most.

While the GAN model generally performs better than both VAE variants, it is worth noting that this advantage may depend on the synthetic data characteristics, which align well with its inductive biases. It achieves substantial gains in IoU (rising from 0.35 to above 0.55), F1-score (approaching 0.75), and recall (peaking above 0.80), while maintaining high precision. The GAN’s edge likely stems from its ability to retain fine spatial detail

and manage overlapping structures—an outcome presumably supported by the U-Net backbone and the residual hierarchy. Still, how much each design choice contributes individually remains to be studied in more depth. Although the GAN does not achieve perfection, its consistent advantage across all metrics strongly suggests that it captures structural and spectral nuances more effectively than the VAE-based models.



**Figure 4.** Validation performance metrics over 30 training epochs for three evaluated models: **(top)** baseline VAE, **(center)** improved VAE, and **(bottom)** proposed GAN. Metrics include Intersection-over-Union (IoU), F1-score, precision, and recall.

#### 4.2. Discussion

The superior performance of our GAN architecture derives from three key architectural innovations that address fundamental challenges in document restoration. Skip connections between encoder and decoder layers enable exceptional texture preservation by maintaining high-frequency graphical features often lost in conventional approaches. Our hybrid loss function combines adversarial and perceptual components to ensure balanced optimization, effectively preventing mode collapse while maintaining training stability. Crucially, the integration of Instance Normalization enables robust spectral generalization across diverse writing systems by adaptively handling script-specific variability in stroke morphology and ink density.

The reconstruction weighting parameter ( $\lambda = 100$ ) was empirically determined to optimally balance character integrity preservation against artifact suppression, particularly for delicate historical manuscripts where excessive smoothing would obliterate critical



paleographic details. This configuration maintains sharp glyph boundaries while effectively eliminating noise patterns characteristic of degraded parchment substrates. Unlike physics-based methods [6], this approach does not require hardware-specific calibration. Compared to encoder–decoder architectures [8], it achieves 34.7% higher IoU through adversarial training.

While the proposed method demonstrates robust performance across most historical documents, certain constraints merit consideration. The approach shows degraded performance when processing extremely faded texts with contrast levels below 5%, where essential stroke information becomes indistinguishable from parchment noise. Additionally, the model's effectiveness depends on careful script balancing during the training phase, as unequal representation of writing systems can bias the learned features toward more frequent scripts. These limitations suggest promising directions for future work in ultra-low-contrast document recovery and unsupervised script adaptation.

## 5. Conclusions

This research establishes a novel adversarial learning framework that fundamentally advances the computational reconstruction of historically significant manuscripts containing erased and overwritten texts, called palimpsests. Our approach demonstrates three paradigm-shifting contributions to cultural heritage informatics. First, we introduce a physically grounded synthetic data generation methodology that accurately models key degradation processes, including ink diffusion governed by Fickian dynamics, parchment–fiber interactions, and spectral superposition effects. This synthetic framework provides a scalable alternative to scarce real-world training data while maintaining biophysical plausibility validated through micro-CT comparisons.

Second, the proposed adversarial architecture achieves a 34.7% improvement in Intersection-over-Union compared to state-of-the-art variational methods, with significant performance gains particularly evident in complex scripts like Syriac (F1: 0.7357 vs. 0.5810) and Gothic minuscules. This performance advantage stems from the hierarchical feature fusion mechanism that preserves stroke-level details through asymmetric skip connections, coupled with the Markovian discriminator's ability to enforce local texture realism across  $70 \times 70$  receptive fields. Crucially, the framework eliminates the dependency on specialized multispectral imaging hardware, making high-fidelity reconstruction accessible through conventional digital photography.

Future research will focus on three scientifically grounded extensions: semi-supervised domain adaptation using limited real manuscript exemplars via maximum mean discrepancy regularization; temporal modeling of degradation processes through neural differential equations capturing time-dependent diffusion; and multi-spectral generalization extending the framework to hyperspectral (300+ channel) reconstruction. These advances will further bridge the gap between computational innovation and practical philological application, providing scholars with previously inaccessible textual evidence while preserving fragile cultural heritage through non-invasive digital analysis.

As part of future work, the authors intend to conduct a comprehensive ablation study to quantitatively evaluate the incremental contribution of each module in our system. This research will involve a series of controlled experiments where individual components are systematically removed or replaced, allowing us to isolate and measure their specific impact on overall performance. The results will offer a deeper understanding of the relative importance and interactions between modules, thereby guiding further refinement and optimization of the architecture.

**Author Contributions:** Conceptualization, J.L.S.; methodology, J.L.S.; validation, J.L.S. and E.F.-P.; formal analysis, J.L.S.; writing—original draft preparation, J.L.S. and E.F.-P.; writing—review and editing, J.L.S. and E.F.-P. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research received no external funding.

**Data Availability Statement:** The original contributions presented in this study are included in the article. Further inquiries can be directed to the corresponding author.

**Conflicts of Interest:** The authors declare no conflicts of interest.

## References

- Emery, D.; Easton, R. Spectral Imaging and Analytical Approaches for Palimpsest Research. *J. Cult. Herit.* **2021**, *48*, 129–138.
- Seales, W.B.; Parker, C.; Segal, M.; Tov, E.; Shor, P.; Porath, Y. From Damage to Discovery: Virtual Unwrapping of Damaged Manuscripts. *IEEE Signal Process. Mag.* **2016**, *33*, 28–37.
- Jampour, M. Revealing Palimpsests with Latent Diffusion Models: A Generative Approach to Image Inpainting and Handwriting Reconstruction. In Proceedings of the 2025 IEEE/CVF Winter Conference on Applications of Computer Vision Workshops (WACVW), Tucson, AZ, USA, 28 February–4 March 2025; pp. 242–249.
- Kingma, D.P.; Welling, M. Auto-Encoding Variational Bayes. *arXiv* **2013**, arXiv:1312.6114.
- Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; Bengio, Y. Generative Adversarial Networks. *arXiv* **2014**, arXiv:1406.2661. [[CrossRef](#)]
- Mitra, A.; Roy, S.; Bhattacharya, U. Multispectral Document Imaging: A Survey. *Comput. Vis. Image Underst.* **2021**, *210*, 103245.
- Perino, M.; Pronti, L.; Moffa, C.; Rosellini, M.; Felic, A.C. New Frontiers in the Digital Restoration of Hidden Texts in Manuscripts: A Review of the Technical Approaches. *Heritage* **2024**, *7*, 683–696. [[CrossRef](#)]
- Chen, J.; Yu, W.; Sun, K.; Li, C.; Wang, J. Document Image Enhancement using Generative Adversarial Networks. *Pattern Recognit. Lett.* **2021**, *152*, 82–88.
- Bhowmik, S. Document Image Binarization. In *Document Layout Analysis*; SpringerBriefs in Computer Science; Springer: Singapore, 2019; pp. 11–30. [[CrossRef](#)]
- Isola, P.; Zhu, J.Y.; Zhou, T.; Efros, A. Image-to-Image Translation with Conditional Adversarial Networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Honolulu, HI, USA, 21–26 July 2017; pp. 1125–1134.
- Ronneberger, O.; Fischer, P.; Brox, T. U-Net: Convolutional Networks for Biomedical Image Segmentation. In Proceedings of the 18th International Conference, Munich, Germany, 5–9 October 2015; pp. 234–241.
- Woo, S.; Park, J.; Lee, J.Y.; Kweon, I.S. CBAM: Convolutional Block Attention Module. In Proceedings of the European Conference on Computer Vision, Munich, Germany, 8–14 September 2018; pp. 3–19.
- Zhang, H.; Goodfellow, I.; Metaxas, D.; Odena, A. Self-Attention Generative Adversarial Networks. In Proceedings of the 36th International Conference on Machine Learning, Long Beach, CA, USA, 10–15 June 2019; pp. 7354–7363.
- Heusel, M.; Ramsauer, H.; Unterthiner, T.; Nessler, B.; Hochreiter, S. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. *Adv. Neural Inf. Process. Syst.* **2017**, *30*, 6629–6640.
- Starynska, A.; Messinger, D.; Kong, Y. Revealing a history: Palimpsest text separation with generative networks. *Int. J. Doc. Anal. Recognit.* **2021**, *24*, 181–195. [[CrossRef](#)]
- Bird, R.B.; Stewart, W.E.; Lightfoot, E.N. *Transport Phenomena*, 2nd ed.; Contains Derivation and Discussion of Fick’s Second Law of Diffusion; John Wiley & Sons: New York, NY, USA, 2002.
- Han, Y.; Hamon, F.P.; Jiang, S.; Durlflosky, L.J. Surrogate model for geological CO storage and its use in hierarchical MCMC history matching. *Adv. Water Resour.* **2024**, *187*, 104678. [[CrossRef](#)]
- Zhang, Y.; Araya-Polo, M.; Mukerji, T. Synthetic Data Generation for Deep Learning-Based Seismic Inversion: From 1D to Complex 2D Models. *Geophysics* **2022**, *87*, R507–R522.
- Karniadakis, G.E.; Kevrekidis, I.G.; Lu, L.; Perdikaris, P.; Wang, S.; Yang, L. Physics-informed machine learning. *Nat. Rev. Phys.* **2021**, *3*, 422–440. [[CrossRef](#)]
- Arridge, S.; de Hoop, M.; Maass, P.; Öktem, O.; Schönlieb, C.; Unser, M. Deep Learning and Inverse Problems. *Snapshots Mod. Math. Oberwolfach* **2019**, *15*. [[CrossRef](#)]
- Zheng, X.; Xu, Z.; Yin, Q.; Bao, Z.; Chen, Z.; Wang, S. A Transformer-Unet Generative Adversarial Network for the Super-Resolution Reconstruction of DEMs. *Remote Sens.* **2024**, *16*, 3676. [[CrossRef](#)]
- Burgess, C.P.; Higgins, I.; Pal, A.; Matthey, L.; Watters, N.; Desjardins, G.; Lerchner, A. Understanding disentangling in  $\beta$ -VAE. *arXiv* **2018**, arXiv:1804.03599.
- Razavi, A.; van den Oord, A.; Vinyals, O. Preventing Posterior Collapse with  $\delta$ -VAEs. In Proceedings of the International Conference on Learning Representations (ICLR), New Orleans, LA, USA, 6–9 May 2019.

24. Zhao, S.; Song, J.; Ermon, S. Towards Deeper Understanding of Variational Autoencoding Models. In Proceedings of the 34th International Conference on Machine Learning, Sydney, Australia, 6–11 August 2017; Volume 70, pp. 3981–3990.
25. Seo, H.j.; Kim, D.; Chung, H.; Lee, S. Handwritten text segmentation via end-to-end learning of convolutional neural network. *Pattern Recognit.* **2020**, *107*, 107473.
26. Ulyanov, D.; Vedaldi, A.; Lempitsky, V. Instance normalization: The missing ingredient for fast stylization. *arXiv* **2016**, arXiv:1607.08022.
27. Srivastava, N.; Hinton, G.; Krizhevsky, A.; Sutskever, I.; Salakhutdinov, R. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *J. Mach. Learn. Res.* **2014**, *15*, 1929–1958.
28. Bowman, S.R.; Vilnis, L.; Vinyals, O.; Dai, A.M.; Jozefowicz, R.; Bengio, S. Generating sentences from a continuous space. In Proceedings of the CoNLL, Beijing, China, 30–31 July 2015.
29. Sudre, C.H.; Li, W.; Vercauteren, T.; Ourselin, S.; Cardoso, M.J. Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations. In Proceedings of the International Conference of Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support, DLMIA ML-CDS 2017, Québec City, QC, Canada, 14 September 2017; pp. 240–248. [\[CrossRef\]](#)
30. Lin, T.Y.; Goyal, P.; Girshick, R.; He, K.; Dollár, P. Focal loss for dense object detection. In Proceedings of the 2017 IEEE International Conference on Computer Vision (ICCV), Venice, Italy, 22–29 October 2017. [\[CrossRef\]](#)
31. Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A.N.; Kaiser, Ł.; Polosukhin, I. Attention is All You Need. In Proceedings of the Advances in Neural Information Processing Systems (NeurIPS), Long Beach, CA, USA, 4–9 December 2017; pp. 5998–6008.
32. Mao, X.; Li, Q.; Xie, H.; Lau, R.; Wang, Z. Least Squares Generative Adversarial Networks. In Proceedings of the IEEE International Conference on Computer Vision, Venice, Italy, 22–29 October 2017.
33. Gulrajani, I.; Ahmed, F.; Arjovsky, M.; Dumoulin, V.; Courville, A. Improved Training of Wasserstein GANs. In Proceedings of the 31st International Conference on Neural Information Processing Systems, Long Beach, CA, USA, 4–9 December 2017.
34. Kingma, D.P.; Ba, J. Adam: A method for stochastic optimization. *arXiv* **2014**, arXiv:1412.6980.
35. Radford, A.; Metz, L.; Chintala, S. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv* **2015**, arXiv:1511.06434.
36. You, Y.; Gitman, I.; Ginsburg, B. Large batch training of convolutional networks. *arXiv* **2017**, arXiv:1708.03888. [\[CrossRef\]](#)
37. He, K.; Zhang, X.; Ren, S.; Sun, J. Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification. In Proceedings of the IEEE International Conference on Computer Vision (ICCV), Washington, DC, USA, 7–13 December 2015; pp. 1026–1034.
38. Wang, Z.; Bovik, A.; Sheikh, H.; Simoncelli, E. Image quality assessment: From error visibility to structural similarity. *IEEE Trans. Image Process.* **2004**, *13*, 600–612. [\[CrossRef\]](#) [\[PubMed\]](#)
39. Rabaev, I.; Litvak, M. Recent advances in text line segmentation and baseline detection in historical document images: A systematic review. *Int. J. Doc. Anal. Recognit.* **2025**. [\[CrossRef\]](#)
40. Powers, D.M. Evaluation: From precision, recall and F-measure to ROC, informedness, markedness and correlation. *J. Mach. Learn. Technol.* **2011**, *2*, 37–63.
41. Otsu, N. A threshold selection method from gray-level histograms. *IEEE Trans. Syst. Man Cybern.* **1979**, *9*, 62–66. [\[CrossRef\]](#)

**Disclaimer/Publisher’s Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.