

VAE-CYCLEGAN: VARIATIONAL LATENT REPRESENTATION FOR UNPAIRED IMAGE-TO-IMAGE TRANSLATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Image-to-image translation plays a central role in computer vision, enabling applications such as style transfer, domain adaptation, and image enhancement. While recent advances have achieved strong paired translation results, learning mappings in unpaired settings remains challenging. In this work, we present a systematic comparison of autoencoder and variational autoencoder (VAE) variants for unpaired image-to-image translation, using paired data solely as a reference baseline. To capture distributional uncertainty, we introduce VAE-CycleGAN, a unified probabilistic framework that integrates variational inference into the CycleGAN architecture. Our method combines adversarial training and cycle-consistency with a VAE’s probabilistic latent space, allowing the model to approximate the true posterior distribution. Further, the architecture achieves a $256\times$ spatial compression, efficiently compressing the input into a compact latent representation. Empirical results across the satellite-to-map benchmark dataset demonstrate that VAE-CycleGAN generates high-quality translated images (FID: 67.75) and achieves superior reconstruction fidelity (MSE: 0.0010, PSNR: 29.85 dB, SSIM: 0.7873) comparable to state-of-the-art deterministic approaches without hyperparameter tuning. For summer-to-winter and label-to-cityscape datasets, VAE-CycleGAN performs comparably with state-of-the-art UNSB at 1 step, and is far superior to UNIT-DDPM at 1000 steps, while the deterministic AE-CycleGAN is comparable to the 5-step UNSB variant.

1 INTRODUCTION

Image-to-image translation, which involves mapping images from a source domain to a target domain while preserving content, has been widely studied in paired (Isola et al., 2017) and unpaired settings (Zhu et al., 2017a). While paired translation assumes aligned image pairs, datasets often lack such correspondence.

Applications like image design and artistic style transfer may have paired data, but scientific fields suffer from limited measurements, and unknown physical processes, making paired real-world or synthetic simulation data infeasible. Some examples include atmospheric remote sensing (satellite to maps, cloud cover, or vegetation), geophysical inversion (images of seismic wave velocity to natural resources), fluids (laminar to turbulent flow, images to velocity fields), and robotic vision (cityscape to labels). All these domains involve information asymmetry, leading to ill-posed translation.

Unpaired translation is thus more challenging. To this end, CycleGAN (Zhu et al., 2017a) introduced cycle-consistency for self-supervision. However, their deterministic mapping fails to capture the inherently multi-modal nature of many ill-posed translation tasks, where a single input corresponds to multiple outputs. Stochastic models have different problems: state-of-the-art diffusion models Ho et al. (2020) produce high fidelity images but have difficulty preserving positional information, with features slightly misaligned or shifted. Visually, these issues are negligible, but for scientific applications, computations such as mass, momentum, energy, area, volume can deteriorate, especially for chaotic simulations.

Meanwhile, nearly all state-of-the-art unpaired diffusion approaches rely either on coarse text-based manipulation or on a “translation module,” typically implemented as an autoencoder. In practice,

054 achieving full posterior sampling requires semantic or latent noise, which in turn necessitates a pre-
055 trained VAE. Thus, any high-performing unpaired translation model is implicitly built upon com-
056 ponents closely aligned with the architecture we propose, even when not presented as such. Rather
057 than posing a strict choice between VAEs and diffusion models, our findings indicate that a strong
058 VAE is essential for building a strong diffusion model. We explored connecting the VAE directly to a
059 diffusion process, following DiffuseVAE, and—even with only a tenth of the training—we observed
060 more than a 6-point FID improvement over the base VAE’s 65 FID.

061 Motivated by these insights and by the inherent multi-modality and positional sensitivity of sci-
062 entific image-to-image translation tasks, we introduce VAE-CycleGAN, a unified framework that
063 integrates the probabilistic latent space of a VAE into the CycleGAN architecture. Our model com-
064 bines adversarial and cycle-consistent training with a structured latent distribution to enable diverse,
065 controllable, and physically faithful translation. We further show that our deterministic variant, AE-
066 CycleGAN, attains state-of-the-art fidelity—surpassing several U-Net and diffusion baselines, while
067 preserving the precise spatial information crucial for scientific applications

068 2 RELATED WORK

069 Generative Adversarial Networks (GANs) (Goodfellow et al., 2014) introduce a transformative
070 adversarial framework where a generator (G) and discriminator (D) are trained simultaneously
071 through a minimax game. This approach produces high-fidelity samples efficiently via backprop-
072 agation, bypassing the need for Markov chains or explicit likelihoods. Subsequent developments
073 improve training stability and sample quality (Radford et al., 2016; Arjovsky et al., 2017). Ad-
074 vancements on conditional GANs (cGANs) by Isola et al. (2017) enable directed generation but
075 require paired data (x, y) . For high-resolution image generation, Wang et al. (2018) incorporates
076 segmentation information, while Park et al. (2019) advances semantic manipulation capabilities.
077 The challenge of multimodality in paired settings is addressed by Zhu et al. (2017b) with Bicycle-
078 GAN, which combines cGANs with variational objectives to generate diverse outputs from single
079 inputs.

080 However, paired data requirements render cGANs unsuitable for ill-posed inverse problems such as
081 tomographic reconstruction or image colorization, where single inputs correspond to multiple valid
082 outputs, as well as for distribution-level translation between unpaired domains like artistic style
083 transfer. To address paired data limitations, unsupervised methods emerge. CycleGAN (Zhu et al.,
084 2017a) and DualGAN (Yi et al., 2017) pioneer cycle-consistency loss with adversarial training for
085 unpaired bi-directional mapping. UNIT (Liu et al., 2017) introduces shared-latent space assump-
086 tions, while StarGAN (Choi et al., 2018; 2020) enables multi-domain translation within a unified
087 framework.

088 Subsequent efforts addressing deterministic limitations include MUNIT (Huang et al., 2018) and
089 DRIT (Lee et al., 2018), which explicitly disentangle style and content codes for multimodal trans-
090 lation, albeit with increased architectural complexity. Meanwhile, Jha et al. (2018) combines a VAE
091 with cycle-consistency, enforcing the reconstruction of an image after its latent factors (like style
092 and content) have been translated across domains. Specifically, the latent space is disentangled into
093 two complementary subspaces via weak supervision with pairwise similarity labels. Alternatively,
094 Augmented CycleGAN (Almahairi et al., 2018) augments the latent space, enabling multimodality
095 by cycling through different augmentations. Here, contrastive learning approaches as in Park et al.
096 (2020) can leverage patch-wise losses between these augmented spaces.

097 Recently, flow-based models, including normalizing flows and diffusion models provide an alterna-
098 tive generative modeling paradigm through quasi-invertible transformations (Rezende & Mohamed,
099 2015; Kingma & Dhariwal, 2018), enabling both precise likelihood estimation and bi-directional
100 latent space manipulation. Diffusion models (Ho et al., 2020; Song et al., 2021; Saharia et al.,
101 2022; Zhang et al., 2023) utilize iterative denoising (stochastic) processes to achieve state-of-the-art
102 performance in high-resolution image synthesis.

103 For unpaired translation, UNIT-DDPM (Sasaki et al., 2021) employs dual denoising diffusion mod-
104 els trained jointly via conditional score-matching. They find improved Fréchet Inception Distance
105 (FID) performance over prior approaches like CycleGAN and VAEs, though later retraining shows
106 similar CycleGAN performance (Kim et al., 2024).

EGSDE (Zhao et al., 2022) models image translation through energy-guided stochastic differential equations (SDEs), offering enhanced training stability and stochasticity not present in CycleGAN. Similarly, Unpaired Image-to-Image Translation via Neural Schrödinger Bridge (NSB) (Sasaki et al., 2021) frames translation as continuous-time stochastic interpolation. Both of these approaches forgo an explicit, interpretable latent space for higher visual fidelity (FID), and generally outperform CycleGAN. As we show, this need not be the case; explicit, stochastic latent space autoencoders can provide near-equivalent performance with better architecture.

In particular, modern VAE architectures have substantially advanced beyond the original formulation (Kingma et al., 2013). Hierarchical VAEs (Vahdat & Kautz, 2020) employ multi-scale latent representations to capture complex data distributions. Vector-quantized VAEs (VQ-VAE) (Van den Oord et al., 2017) introduce discrete latent representations through a codebook mechanism, effectively circumventing the posterior collapse problem common in continuous VAEs when paired with powerful decoders. This approach replaces the continuous latent space with discrete codes learned via vector quantization, creating a robust framework for high-quality image, video, and speech generation. The subsequent VQ-VAE-2 (Razavi et al., 2019) enhances this foundation through a multi-scale hierarchical architecture, employing powerful autoregressive priors over the latent codes to generate high-fidelity, diverse samples that rival state-of-the-art GANs, while maintaining the training stability and diversity advantages of VAE-based approaches. The NVAE architecture (Vahdat & Kautz, 2020) uses depth-wise separable convolutions to also demonstrate that hierarchical VAEs can compete with other state-of-the-art generative models; the Very Deep VAE (Child, 2021) shows that extremely deep hierarchical variational models can rival autoregressive approaches in sample quality.

Finally, the VAE-GAN paradigm (Larsen et al., 2016) combines variational autoencoders’ latent space learning with GANs’ adversarial training, using the discriminator for perceptually-aware reconstruction losses rather than pixel-wise metrics. This hybrid approach demonstrates enhanced generalization and output fidelity (Yan et al., 2025; Denton & Fergus, 2019), with the VAE learning meaningful latent representations while the adversarial component ensures distributional alignment.

Our proposed VAE-CycleGAN builds upon these advances, incorporating cycle-consistent adversarial training for bi-directional unpaired translation while leveraging VAE-based generators to introduce multimodality and structured latent spaces. This distinguishes our work from standard VAE-GANs (Larsen et al., 2016; Yan et al., 2025) and addresses CycleGAN’s determinism through intrinsic stochasticity rather than external auxiliary variables (Almahairi et al., 2018). While we do not explicitly include diffusion models, our VAE-based model allows for easy integration with latent diffusion models as in Zhang et al. (2023) or refinements like Pandey et al. (2022).

Concurrent research by Sharma et al. (2025) explored medical image translation with unpaired data using a similar combination of a VAE and a CycleGAN. Their architecture employs two GANs where the generator module of only one GAN incorporates a VAE neural network. In contrast, our network integrates a VAE into each of the two GANs, fully leveraging the potential for variability in the translated images. Furthermore, our paper provides a comprehensive comparison of different AE (autoencoder) and VAE variants in terms of reconstruction, translation, and diversity of the translated samples, with both perception (visual quality) and distortion (pointwise accuracy) metrics.

3 PROBLEM FORMULATION

In consistent image translation, we expect a canonical isomorphism between two visual domains X, Y . In practice, all possible such mappings are lossy and thus non-invertible; we thus seek the maximally structure-preserving map (or pair of homomorphisms with minimal kernel). These maps are not unique; given an image $x_i \in X$, we can define the possible outcomes $\hat{y} \sim p_{\hat{Y}|X}$ with a posterior distribution, as in classical ill-posed inverse problems.

As in CycleGAN (Zhu et al., 2017a), we indirectly enforce maximal structure preservation through a *cycle consistency loss*, for generators $G : X \rightarrow Y, F : Y \rightarrow X$.

$$\mathcal{L}_{\text{cycle}} = \mathbb{E}_x [|x - F(G(x))|_1] + \mathbb{E}_y [|y - G(F(y))|_1] \quad (1)$$

Since this is a distortion metric, our model follows the perception-distortion tradeoff (Blau & Michaeli, 2018): an estimator cannot simultaneously achieve optimal accuracy (minimal distortion) and optimal perception (distributional or visual quality). For an allowable distortion level D ,

perception metric (f-divergence, e.g. Kullback-Leibler (KL)) d , and distortion metric (e.g. L1 loss) Δ , the optimal model satisfies Eqn. 2:

$$\hat{p}_{\hat{y}|x} = \operatorname{argmin}_{p_{\hat{y}|x}} d(p_y, p_{\hat{y}}) \quad \text{s.t.} \quad \mathbb{E}_{x,y \sim p_{\text{data}}} \mathbb{E}_{\hat{y} \sim p_{\hat{y}|x}} [\Delta(y, \hat{y})] \leq D \quad (2)$$

Low distortion estimators (e.g., standard VAEs, autoencoders) minimize $\mathbb{E}[\Delta(y, \hat{y})]$ but ignore the distribution $p_{\hat{y}}$, resulting in blurry, perceptually unrealistic outputs (converging to the pixel-wise maximum a-posteriori solution). In contrast, high perception estimators (e.g., GANs) learn and minimize $d(p_y, p_{\hat{y}})$ and produce sharp, realistic samples but provide no guarantee that a generated sample \hat{y} is a faithful reconstruction of a specific input x . An optimal inversion requires sampling from the full posterior $p(y | x)$, which defines the Pareto frontier in Eqn. 2. VAE-CycleGAN is designed to learn this posterior distribution, enabling both perceptually realistic and distortion-faithful reconstructions. We use for perception metric $d(p_y, p_{\hat{y}})$ an *adversarial χ^2 divergence* between generated images and target domains, with discriminators $D_X : X \rightarrow \mathbb{P}_X, D_Y : Y \rightarrow \mathbb{P}_Y$.

$$\mathcal{L}_{\text{GAN}}^{X \rightarrow Y} = \mathbb{E}_y [D_Y(y)^2] + \mathbb{E}_x [(1 - D_Y(G(x)))^2] \quad (3)$$

$$\mathcal{L}_{\text{GAN}}^{Y \rightarrow X} = \mathbb{E}_x [D_X(x)^2] + \mathbb{E}_y [(1 - D_X(F(y)))^2] \quad (4)$$

Interestingly, at no point have we required paired data; by quantifying both perception and distortion the model is now fully specified, even in the unpaired setting (where we have datasets: $\{(x_i \in X) \sim p_{\text{data}}(x)\}_{i=1}^N$ and $\{(y_j \in Y) \sim p_{\text{data}}(y)\}_{j=1}^M$). Please see Appendix 7.1 for a proof sketch of model convergence to a natural map, though not necessarily a human-preferred, canonical map.

We finish the objectives with regularizers. For fast posterior sampling, a β -VAE *loss* regularizes the latent space to a normal distribution with a KL divergence d_{KL} (Higgins et al., 2017), (Weng, 2018). For latent variable z , ϕ and θ parameterize the encoder $q_\phi(z|\cdot)$ and decoder $p_\theta(\cdot|z)$ respectively. For stability, we add an *identity loss* regularizer to preserve content during adversarial training.

$$\mathcal{L}_{\text{VAE}}^X = \mathbb{E}_{z \sim q_\phi(z|x)} d_{\text{KL}}(q_\phi(z|x) \| p(z) \triangleq \mathcal{N}(0, \mathbb{I})) \quad (5)$$

$$\mathcal{L}_{\text{VAE}}^Y = \mathbb{E}_{z \sim q_\phi(z|y)} d_{\text{KL}}(q_\phi(z|y) \| p(z) \triangleq \mathcal{N}(0, \mathbb{I})) \quad (6)$$

$$\mathcal{L}_{\text{identity}} = \mathbb{E}_x [|G(x) - x|_1] + \mathbb{E}_y [|F(y) - y|_1] \quad (7)$$

The complete training objectives are then as in Eqns. 8, 9, where $\lambda_{\text{GAN}}, \lambda_{\text{cycle}}, \lambda_{\text{id}}, \lambda_{\text{kl}} \ll 1$ are weighting coefficients.

$$\text{AE-CycleGAN: } \mathcal{L}_{\text{Total}} = \lambda_{\text{GAN}} \mathcal{L}_{\text{GAN}} + \lambda_{\text{cycle}} \mathcal{L}_{\text{cycle}} + \lambda_{\text{id}} \mathcal{L}_{\text{identity}} \quad (8)$$

$$\text{VAE-CycleGAN: } \mathcal{L}_{\text{Total}} = \lambda_{\text{kl}} \mathcal{L}_{\text{VAE}} + \lambda_{\text{GAN}} \mathcal{L}_{\text{GAN}} + \lambda_{\text{cycle}} \mathcal{L}_{\text{cycle}} + \lambda_{\text{id}} \mathcal{L}_{\text{identity}} \quad (9)$$

3.1 NETWORKS

We train the following variants in Table 1 to fully ablate these objectives, on an asymmetric translation benchmark between high-resolution satellite (aerial) photos X and simplified map-like representations Y , with comprehensive architecture as in Figure 1. We closely follow the CycleGAN architecture, but remove the U-Net skip connections and increase compression by 12x. Upsampling and downsampling instead use pixel shuffle-unshuffle as in Chen et al. (2024). Please see Appendix 7.2, 7.3 for exact implementation details, respectively.

Table 1: Architecture variants

	Model	Type	Objectives
Paired Dataset	AE	Deterministic	Translation
	Cycle AE	Deterministic	Translation, Cycle consistency
	AE-GAN	Deterministic	Translation, Adversarial, Identity
	VAE	Stochastic	Translation, KL-Divergence
	Cycle-VAE	Stochastic	Translation, Cycle consistency, KL-Divergence
	VAE-GAN	Stochastic	Translation, Adversarial, Identity, KL-Divergence
Unpaired	Cycle AE	Deterministic	Cycle consistency
	AE-CycleGAN	Deterministic	Adversarial, Identity, Cycle consistency
	Cycle-VAE	Stochastic	Cycle consistency, KL-Divergence
	VAE-CycleGAN	Stochastic	Adversarial, Identity, Cycle consistency, KL-Divergence

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

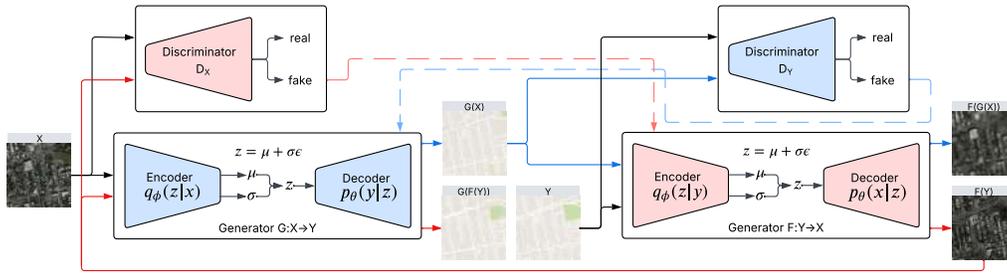


Figure 1: VAE-CycleGAN for unpaired translation. Generators G and F map $X \rightarrow Y$ and $Y \rightarrow X$, respectively, discriminators D_X and D_Y enforce domain alignment, and reconstruction losses $F(G(X)) \approx X$ & $G(F(Y)) \approx Y$ enforce cycle-consistency.

4 RESULTS

We first compare the performance of paired dataset AE and VAE variants, then extend to unpaired settings. We then return to the role of adversarial loss and cycle-consistency in the stochastic (inverse-problem) case, by introducing a Gaussian prior in the latent space.

To evaluate performance, we use distortion-based metrics MSE, PSNR, the perception-oriented Structural Similarity Index (SSIM), and perception-based metrics including Fréchet Inception Distance (FID) and Kernel Inception Distance (KID), to measure distributional distance between real and generated images. We report KID scores since our datasets are relatively small with only about 1K training and testing samples.

Note that all our AE/VAE variants are quite stable and rarely collapses; we see almost no mode collapse provided that the loss weights are not extremely unbalanced. As such, we do not / have not finetuned the loss weights and other hyperparameters beyond scaling the cycle consistency to adequately preserve information.

For brevity, we present an examples of VAE-CycleGAN’s translation and reconstruction (unpaired) in Figures 2, 3. Examples of all AE/VAE variants (listed in Table 1) are presented in Appendix 7.4.

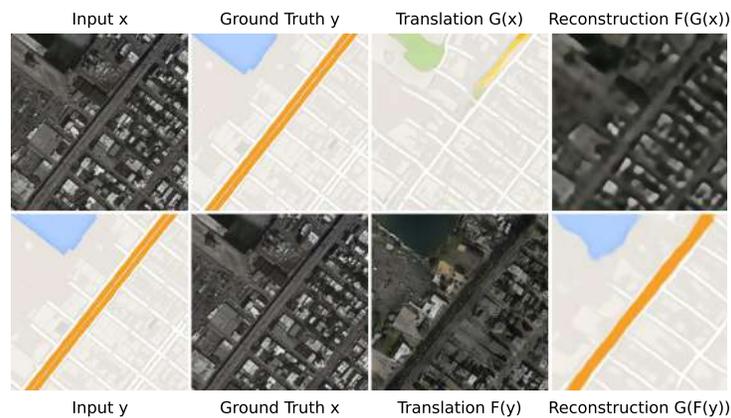


Figure 2: Satellite images \leftrightarrow Maps: VAE-CycleGAN translated and reconstructed outputs.

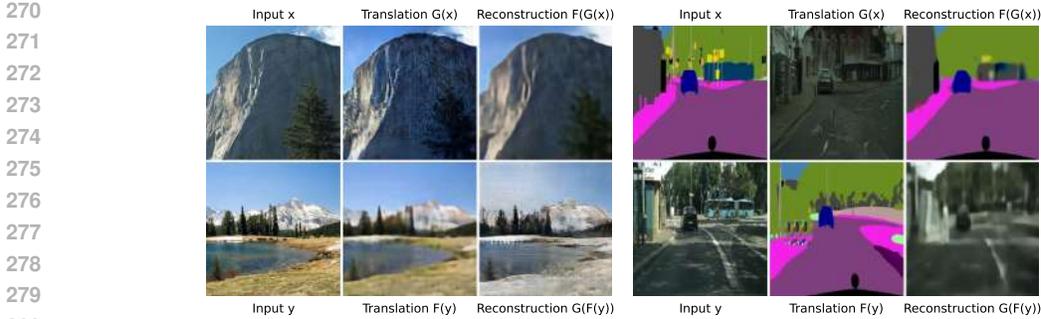


Figure 3: Summer \leftrightarrow Winter, Labels \leftrightarrow Cityscape: VAE-CycleGAN translated and reconstructed outputs. Translations are sharp with the exception of the Winter to Summer case, as the input y could be easily classified as either summer or winter.

Tables 2 and 3 display the average translation and reconstruction errors for the Satellite2Map dataset.

Table 2: Model ablation across translation tasks. KID is scaled by 100.

Model	$G : X \rightarrow Y$ (aerial \rightarrow map)					$F : Y \rightarrow X$ (map \rightarrow aerial)					
	MSE \downarrow	PSNR \uparrow	SSIM \uparrow	FID \downarrow	KID $\mu \pm \sigma \downarrow$	MSE \downarrow	PSNR \uparrow	SSIM \uparrow	FID \downarrow	KID $\mu \pm \sigma \downarrow$	
Paired	AE	0.0018	27.37	0.7636	325.56	35.30 \pm 3.74	0.0256	14.08	0.2489	253.05	23.31 \pm 1.96
	Cycle AE	0.0018	27.46	0.7650	273.17	27.29 \pm 2.41	0.0251	14.13	0.2553	241.77	21.75 \pm 1.89
	AE-GAN	0.0031	25.10	0.6684	63.83	3.16 \pm 0.66	0.0349	14.29	0.2079	52.33	1.75 \pm 0.55
	VAE	0.0024	26.03	0.6753	358.70	35.47 \pm 3.10	0.0261	13.50	0.1894	304.34	28.58 \pm 2.09
	Cycle VAE	0.0022	26.58	0.7117	325.89	33.69 \pm 2.36	0.0255	13.25	0.2080	259.83	24.03 \pm 1.89
	VAE-GAN	0.0033	24.85	0.6272	83.18	5.10 \pm 0.91	0.0354	14.17	0.1705	64.45	3.01 \pm 0.52
Unpaired	Cycle AE	0.0872	10.59	0.0318	409.27	50.77 \pm 1.63	0.0410	13.69	0.0924	329.51	35.52 \pm 1.75
	Cycle VAE	0.0113	19.47	0.3128	419.51	53.93 \pm 1.94	0.0614	7.90	0.0885	475.98	59.31 \pm 2.74
	AE-CycleGAN	0.0050	22.97	0.6737	70.08	4.60 \pm 1.31	0.0389	13.62	0.1641	52.53	1.70 \pm 0.48
	VAE-CycleGAN	0.0056	22.50	0.5793	90.87	5.94 \pm 0.92	0.0443	13.36	0.0965	69.25	3.78 \pm 0.63

Table 3: Model ablation across reconstruction tasks. KID is scaled by 100.

Model	$G(F(y))$: Map Reconstruction					$F(G(x))$: Aerial Reconstruction					
	MSE \downarrow	PSNR \uparrow	SSIM \uparrow	FID \downarrow	KID $\mu \pm \sigma \downarrow$	MSE \downarrow	PSNR \uparrow	SSIM \uparrow	FID \downarrow	KID $\mu \pm \sigma \downarrow$	
Paired	AE	0.0017	27.77	0.7718	315.15	33.47 \pm 3.24	0.0287	13.95	0.1992	288.35	28.07 \pm 2.08
	Cycle AE	0.0003	35.91	0.9242	83.08	6.04 \pm 0.66	0.0059	21.53	0.6661	175.29	16.21 \pm 1.64
	AE-GAN	0.0050	22.99	0.6969	74.76	4.25 \pm 0.82	0.0392	13.90	0.1481	70.75	3.47 \pm 0.67
	VAE	0.0049	22.63	0.6717	367.13	38.95 \pm 2.50	0.0340	11.13	0.1344	347.17	36.42 \pm 1.88
	Cycle VAE	0.0007	31.81	0.8450	170.27	15.46 \pm 1.54	0.0145	17.03	0.3330	266.76	25.87 \pm 1.73
	VAE-GAN	0.0060	22.22	0.6172	104.49	8.00 \pm 1.07	0.0430	13.35	0.0954	84.63	4.96 \pm 1.01
Unpaired	Cycle AE	0.0002	36.59	0.9317	80.05	6.08 \pm 0.75	0.0063	21.12	0.6451	180.12	16.62 \pm 1.88
	Cycle VAE	0.0006	32.44	0.8468	219.37	21.30 \pm 2.82	0.0146	17.31	0.3151	281.64	27.61 \pm 1.68
	AE-CycleGAN	0.0005	32.99	0.8760	106.63	8.44 \pm 0.92	0.0096	19.01	0.5069	200.11	18.33 \pm 1.81
	VAE-CycleGAN	0.0011	29.67	0.7804	241.78	24.09 \pm 2.59	0.0175	16.10	0.2779	271.72	27.03 \pm 2.00

Table 4: Finetuned translation tasks. The AE-CycleGAN is nearly state-of-the-art.

Model	FID \downarrow	Steps	Label2Cityscape		Summer2Winter		Satellite2Map	
			\rightarrow	\leftarrow	\rightarrow	\leftarrow	\rightarrow	\leftarrow
UNIT-DDPM Sasaki et al. (2021)		1000	113.70	109.98	113.70	116.23	193.06	
MUNIT Huang et al. (2018)		1	91.4	115.4	-	-	181.7	
CycleGAN Zhu et al. (2017a)		1	76.3	84.9	-	-	54.6	
CycleDiff Zou et al. (2025)		100	45.1	72.7	-	-	53.2	
UNSB Kim et al. (2024)		1	\approx 74	\approx 75	-	-	\approx 71	
UNSB Kim et al. (2024)		5	<u>53.2</u>	<u>73.9</u>	-	-	47.6	
AE-CycleGAN		1	-	-	-	71.1	49.4	
VAE-CycleGAN (lower λ_{kl})		1	-	-	-	<u>89.8</u>	57.3	
VAE-CycleGAN		1	67.9	86.6	72.4	98.4	67.7	

The asymmetric difficulty of the aerial and map datasets is immediately obvious, with any aerial task generally worse in distortion metrics (high MSE, low PSNR, SSIM) but easier to capture in-distribution (low FID, KID). Regardless, optimal translation and reconstruction models closely follow the perception-distortion trade-off (Blau & Michaeli, 2018). In both paired translation and reconstruction, Cycle-AE minimizes distortion via a cycle-consistency loss (lowest MSE, high PSNR, SSIM), while AE-GAN optimizes for perception with the adversarial loss (lowest FID, KID). Interestingly, in unpaired translation, the adversarial loss becomes crucial for domain alignment. Consequently, AE-CycleGAN achieves the best overall performance, superior in both distortion and perception metrics. Similarly, Cycle AE (with only cycle-consistency) performs best in all reconstruction metrics, as cycle-consistency is crucial for information preservation.

We find AE-CycleGAN the best overall performing (deterministic) model. No metric in Table 2 directly quantifies the conditional posterior (distribution of cycle-consistent translations), since that is intractable over the large 256x256x3 image space, so stochastic models are expected to perform worse on this benchmark. FID and KID only quantify the marginal posterior (any conditioned translation regardless of consistency). Nevertheless, unlike the other stochastic VAE models, VAE-CycleGAN performs closely behind AE-CycleGAN in translation. We expand on this with qualitative measurements in section 5.

Against other state-of-the-art U-Net and diffusion models, AE-CycleGAN achieves comparable fidelity. We outperform the CycleGAN U-Net (Zhu et al., 2017a) at 12x (greater) compression by simply modernizing the autoencoder portion: we remove skip connections and switch to pixel shuffle layers. We also outperform MUNIT (U-Net), UNIT-DDPM and CycleDiff (U-Net with multiple steps). Finally, we outperform the state-of-the-art UNSB (Kim et al., 2024) at 1-step and are comparable with their 5-step model. No comparisons are made against text-conditioned diffusion models; those models are neither quantitative enough for scientific applications, nor are they pixel-level cycle consistent. Also, the VAE-CycleGAN cannot be directly compared with these diffusion models, as state-of-the-art diffusion models are built on top of VAEs; by design, diffusion models need not approximate the complete posterior given limited training data (which is why conditional models and posterior sampling were developed). In some cases, the models actually memorize the training data (Somepalli et al., 2023).

We now visualize translations from aerial photos to maps, $x \rightarrow G(x)$ in Figure 4 and the aerial image reconstructions, $x \approx F(G(x))$ in Figure 5 for additional qualitative comparison.

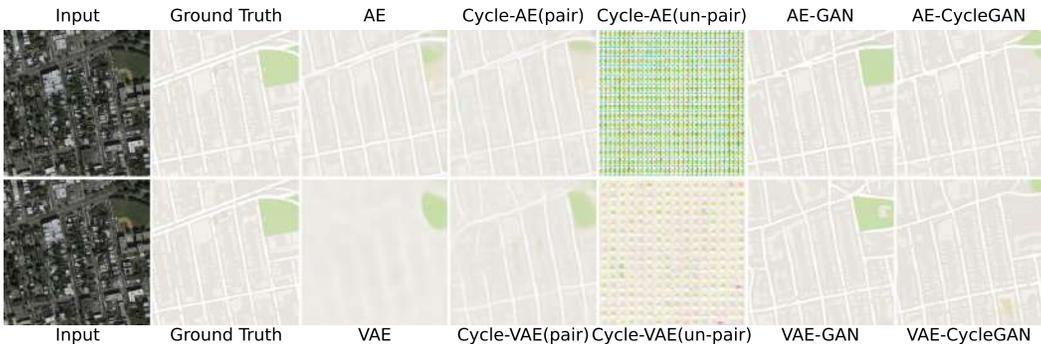
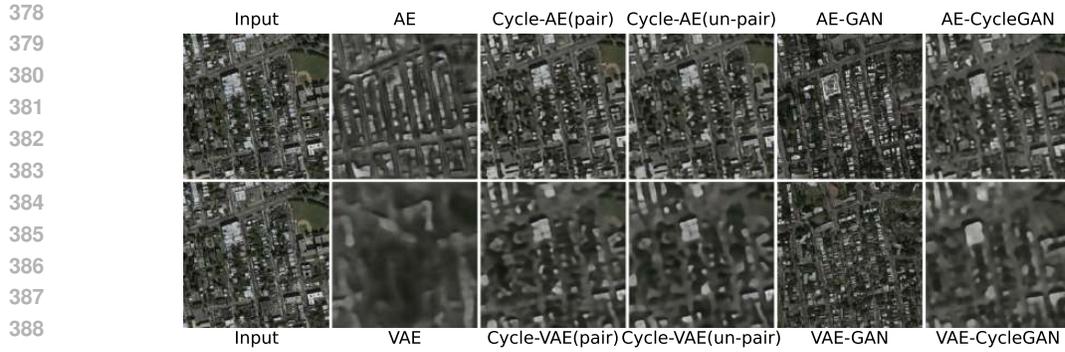


Figure 4: Left to right: input x , target y , map translations $G(x)$. Top to bottom: AE/VAE variants.

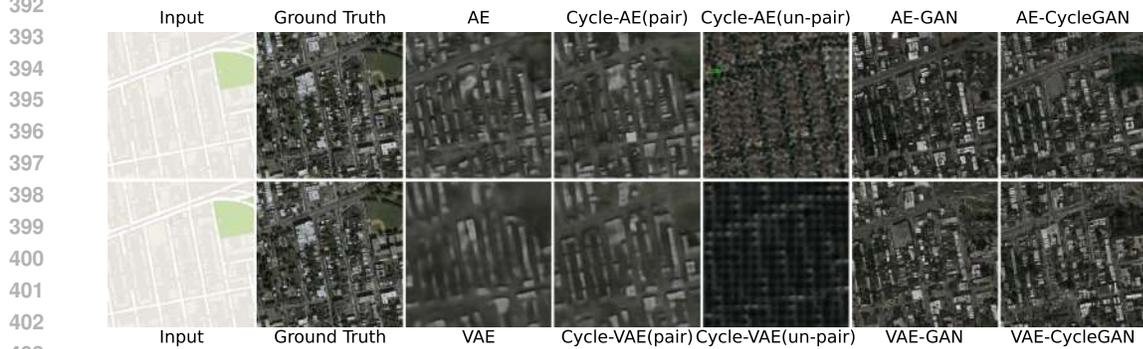
Per Figure 4, maps translated by the Cycle-AE and Cycle-VAE models in an unpaired data setting show no structural similarity to the map domain, as there is neither an adversarial constraint nor a paired dataset to ensure domain alignment. However, spatial information is preserved and allows for the complete reconstruction of the original aerial input, as shown in Figure 5. We notice a similar pattern for aerial photos translation, $y \rightarrow F(y)$ in Figure 6 and map reconstruction, $y \approx G(F(y))$ in Figure 7.

Visually, VAE-CycleGAN and AE-CycleGAN both produce the highest quality translations, with the VAE-CycleGAN enabling distributional sampling at a moderate hit to reconstruction quality, as expected from the perception-distortion tradeoff Blau & Michaeli (2018).



390 Figure 5: Left to right: input x , aerial reconstructions $F(G(x))$. Top to bottom: AE/VAE variants.

391



404 Figure 6: Left to right: input y , target x , aerial translations $F(y)$. Top to bottom: AE/VAE variants.

405

406



418 Figure 7: Left to right: input y and map reconstructions $G(F(y))$. Top to bottom: AE/VAE variants.

419

420

421

422

423 5 REALIZATIONS

424

425

426

427

428

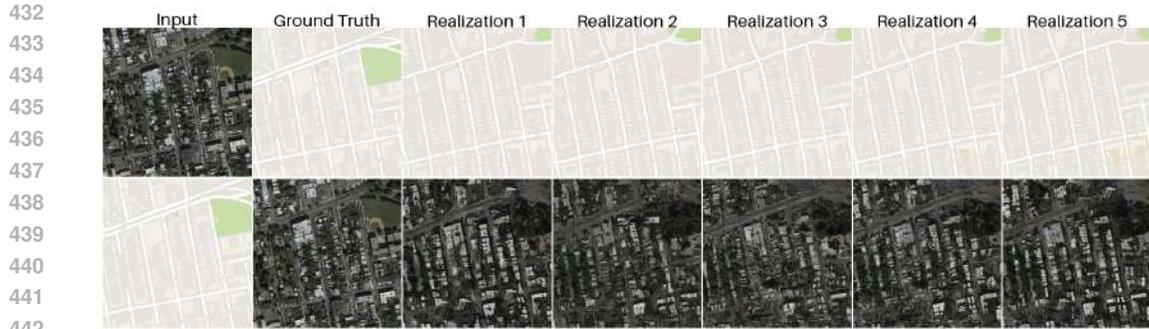
429

430

431

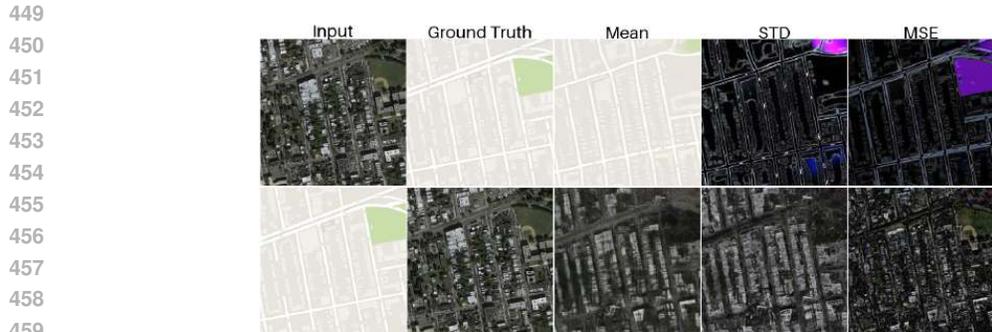
Realizations, ensemble mean, standard deviation, and ensemble error together form a qualitative measure of the intractable conditional performance mentioned in section 4. For brevity, we now visualize VAE-CycleGAN realizations and ensemble statistics in Figures 8, 9. Please find expanded figures and metrics in Appendix 7.5 and 7.6; figures 54 & 53 and 51 & 52 show translations/reconstructions for the cityscape-to-label dataset and summer-to-winter dataset.

As expected, we find a tradeoff in feature sharpness. In aerial output from Figure 8, we find homogeneous, low-information areas (such as roads and open fields) lose sharpness, while sharp edges and textured areas (such as buildings) gain sharpness. Similar behavior is visible in the ensemble mean (Figure 9). Cycle-consistency is satisfied on the low-information areas (even if the translation



443 Figure 8: VAE-CycleGAN sample realizations. Row 1: aerial \rightarrow maps, row 2: maps \rightarrow aerial

444
445
446 is missing information, the details are trivially reconstructed), so it follows that the adversarial loss
447 is poorly aligned with human perception.



460 Figure 9: Mean, scaled STD (aerial $3\times$, maps $20\times$), and MSE across 30 realizations. The standard
461 deviation and MSE are reported in true color, with purple color in standard deviation demonstrating
462 that the generator is confident in green value but highly variable in red and blue. Such behavior is
463 especially prominent in heavily treed or grassy areas (green areas) in satellite images; as expected,
464 the generators learn green value well due to the imbalanced dataset (trees and grass dominate the
465 satellite images).

466
467 As no VAE variant can outperform the best deterministic autoencoder, from the perception-distortion
468 tradeoff, we instead follow standard Bayesian modeling skill. Skill is defined as the correlation of
469 ensemble mean error and standard deviation, and we find that across realizations, the model does
470 produce diverse outputs for regions of uncertainty. In satellite to map translation, uncertainty is
471 highest when translating low-detail ground features (like grass or empty fields) as there is simply
472 not enough information for accurate translation. In map to satellite translation, high-detail textures
473 (like buildings) naturally have highest uncertainty (as an ill-posed inverse problem). Across an
474 ensemble of 30 samples, the map translation MSE is $\approx 6.4 \times 10^{-3}$ while the aerial translation MSE
475 is $\approx 35.1 \times 10^{-3}$ for VAE-CycleGAN, due to the greater ill-posedness of the latter.

476 For all VAE models, Figures 10 and 11 display ensemble summaries. For visibility, map standard
477 deviation images for all networks were scaled by a factor of 20 (excepting VAE-CycleGAN at $25\times$)
478 and then clamped to the 0-1 range for display. Similarly, aerial standard deviation images were
479 scaled by $10\times$ for VAE and Cycle VAE variants, by $3\times$ for GAN based VAEs, and also clamped.

480 In Figures 10 and 11, we observe that the ensemble means for the GAN-based VAEs are much
481 sharper than the VAE and paired Cycle-VAE, despite the paired data. Conversely, the unpaired
482 Cycle-VAE completely lacks domain alignment, without the adversarial loss or paired data. For
483 nearly all models except VAE and the unpaired Cycle-VAE, standard deviation is again a good
484 estimator of ensemble MSE. Lastly, note VAE-CycleGAN shows some positional distortion in the
485 map due to the limited domain alignment possible in the unpaired setting; it is the only performant
unpaired stochastic VAE model.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500

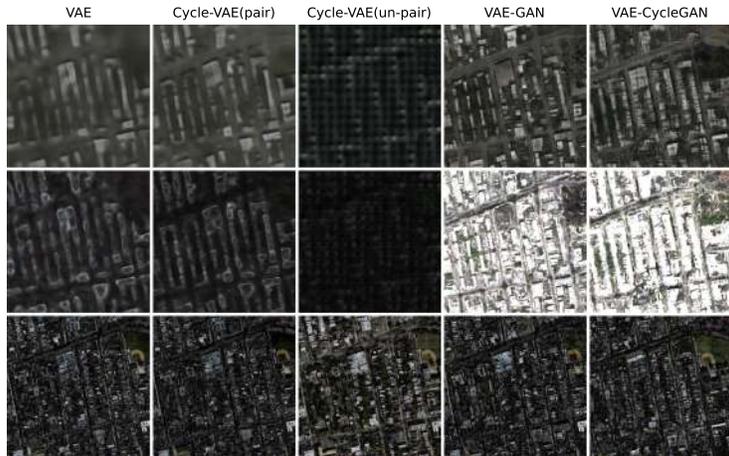


Figure 10: Rows: mean, STD (scaled for visibility), and MSE of 30 aerial realizations (VAE variants).

501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519

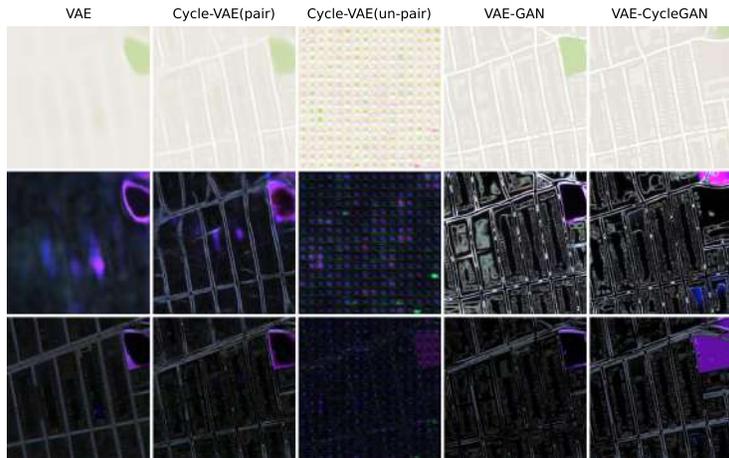


Figure 11: Rows: mean, STD (scaled for visibility), and MSE of 30 map realizations (VAE variants).

6 CONCLUSION

The AE-CycleGAN and VAE-CycleGAN models perform competitively both on distortion and perception metrics during translation. VAE-CycleGAN further maintains output diversity, at an acceptable fidelity tradeoff. The variational latent space allows future adaptation with multi-step methods such as diffusion models, enabling fine-grained or prompt-based control.

As our framework applies variational inference and modern autoencoders to the CycleGAN architecture, extensions to general class of unpaired, bidirectional ill-posed inverse problems are straightforward. Possible applications beyond natural images could include molecular design and medical image synthesis (CT-to-MRI synthesis, PET-to-CT conversion, etc.) from 3D or other non-image data.

ACKNOWLEDGMENTS

The authors acknowledge the use of AI tools for visualization code and language refinement for clarity. All AI-generated content was critically reviewed, adapted, and validated to ensure scientific accuracy, with the authors maintaining full responsibility for the research and intellectual content.

520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

REFERENCES

- Amjad Almahairi, Sai Rajeswar, Alessandro Sordani, Philip Bachman, and Aaron Courville. Augmented cycleGAN: Learning many-to-many mappings from unpaired data. In *International Conference on Machine Learning*, pp. 195–204, Jul 2018.
- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International Conference on Machine Learning*, pp. 214–223, 2017.
- Yochai Blau and Tomer Michaeli. The Perception-Distortion Tradeoff. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6228–6237, June 2018. doi: 10.1109/CVPR.2018.00652. URL <http://arxiv.org/abs/1711.06077>. arXiv:1711.06077 [cs].
- Junyu Chen, Han Cai, Junsong Chen, Enze Xie, Shang Yang, Haotian Tang, Muyang Li, Yao Lu, and Song Han. Deep compression autoencoder for efficient high-resolution diffusion models. *arXiv preprint arXiv:2410.10733*, 2024.
- Rewon Child. Very deep vae generalize autoregressive models and can outperform them on images. In *International Conference on Learning Representations*, 2021.
- Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8789–8797, 2018.
- Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8188–8197, 2020.
- Emily Denton and Rob Fergus. Stochastic video generation with a learned prior. In *International Conference on Machine Learning*, pp. 1574–1583, 2019.
- Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.
- Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *International conference on learning representations*, 2017.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, volume 33, pp. 6840–6851, 2020.
- Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 172–189, 2018.
- Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1125–1134, 2017.
- Ananya Jha, Saket Anand, Maneesh Singh, and VSR Veeravasaru. *Disentangling Factors of Variation with Cycle-Consistent Variational Auto-encoders: 15th European Conference, Munich, Germany, September 8–14, 2018, Proceedings, Part III*, pp. 829–845. 09 2018. ISBN 978-3-030-01218-2. doi: 10.1007/978-3-030-01219-9_49.
- Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual Losses for Real-Time Style Transfer and Super-Resolution. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling (eds.), *Computer Vision – ECCV 2016*, pp. 694–711, Cham, 2016. Springer International Publishing. ISBN 978-3-319-46475-6.
- Beomsu Kim, Gihyun Kwon, Kwanyoung Kim, and Jong Chul Ye. Unpaired image-to-image translation via neural schrödinger bridge, 2024. URL <https://arxiv.org/abs/2305.15086>.

- 594 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua
595 Bengio and Yann LeCun (eds.), *ICLR (Poster)*, 2015. URL [http://dblp.uni-trier.de/
596 db/conf/iclr/iclr2015.html#KingmaB14](http://dblp.uni-trier.de/db/conf/iclr/iclr2015.html#KingmaB14).
- 597 Diederik P. Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions.
598 In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- 600 Diederik P Kingma, Max Welling, et al. Auto-encoding variational bayes, 2013.
- 601 Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther. Au-
602 toencoding beyond pixels using a learned similarity metric, February 2016. URL [http://
603 //arxiv.org/abs/1512.09300](http://arxiv.org/abs/1512.09300). arXiv:1512.09300 [cs].
- 605 Hsin-Ying Lee, Hung-Yu Tseng, Jia-Bin Huang, Maneesh Singh, and Ming-Hsuan Yang. Diverse
606 image-to-image translation via disentangled representations. In *Proceedings of the European
607 Conference on Computer Vision (ECCV)*, pp. 35–51, 2018.
- 609 Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks.
610 In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- 611 Kushagra Pandey, Avideep Mukherjee, Piyush Rai, and Abhishek Kumar. Diffusevae: Effi-
612 cient, controllable and high-fidelity generation from low-dimensional latents. *arXiv preprint
613 arXiv:2201.00308*, 2022.
- 615 Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with
616 spatially-adaptive normalization. In *Proceedings of the IEEE Conference on Computer Vision and
617 Pattern Recognition*, pp. 2337–2346, 2019.
- 618 Taesung Park, Alexei A. Efros, Richard Zhang, and Jun-Yan Zhu. Contrastive learning for unpaired
619 image-to-image translation. In *European Conference on Computer Vision*, pp. 319–345, 2020.
- 621 Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep
622 convolutional generative adversarial networks. In *International Conference on Learning Repre-
623 sentations*, 2016.
- 624 Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with
625 vq-vae-2. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- 627 Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. In
628 *International Conference on Machine Learning*, pp. 1530–1538, 2015.
- 629 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kam-
630 yar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Palette:
631 Image-to-image diffusion models. In *ACM SIGGRAPH 2022 Conference Proceedings*, pp. 1–10,
632 2022.
- 634 Hiroshi Sasaki, Chris G. Willcocks, and Toby P. Breckon. Unit-ddpm: Unpaired image transla-
635 tion with denoising diffusion probabilistic models, 2021. URL [https://arxiv.org/abs/
636 2104.05358](https://arxiv.org/abs/2104.05358).
- 637 Yashvi Sharma, Shikha Diya, and Najme Zehra Naqvi. Images-variational autoencoder cyclegan.
638 *Proceedings of Data Analytics and Management: ICDAM 2024, Volume 3*, 1299:329, 2025.
- 640 Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion
641 art or digital forgery? investigating data replication in diffusion models. In *Proceedings of the
642 IEEE/CVF conference on computer vision and pattern recognition*, pp. 6048–6058, 2023.
- 643 Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben
644 Poole. Score-based generative modeling through stochastic differential equations. In *Interna-
645 tional Conference on Learning Representations*, 2021.
- 647 Arash Vahdat and Jan Kautz. Nvae: A deep hierarchical variational autoencoder. In *Advances in
Neural Information Processing Systems*, volume 33, pp. 19667–19679, 2020.

- 648 Aaron Van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learn-
649 ing. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- 650
- 651 Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-
652 resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of*
653 *the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8798–8807, 2018.
- 654 Lilian Weng. From autoencoder to beta-vae. *lilianweng.github.io*, 2018. URL [https://](https://lilianweng.github.io/posts/2018-08-12-vae/)
655 lilianweng.github.io/posts/2018-08-12-vae/.
- 656
- 657 Jian'en Yan, Haihui Huang, Kairan Yang, Haiyan Xu, and Yanling Li. Synthetic data for enhanced
658 privacy: A vae-gan approach against membership inference attacks. *Knowledge-Based Systems*,
659 309:112899, 2025.
- 660
- 661 Zili Yi, Hao Zhang, Ping Tan, and Minglun Gong. Dualgan: Unsupervised dual learning for image-
662 to-image translation. In *Proceedings of the IEEE International Conference on Computer Vision*,
663 pp. 2849–2857, 2017.
- 664 Lvmin Zhang, Yi Zhang, Jiawei Zhang, Yang Liu, Yifan Huang, Chunyu Wang, Fang Zhao, and
665 Hang Zhou. Diffusion-4k: High-fidelity diffusion model for 4k image generation. In *Proceedings*
666 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12345–12354,
667 2023.
- 668 Min Zhao, Fan Bao, Chongxuan Li, and Jun Zhu. Egsde: Unpaired image-to-image translation
669 via energy-guided stochastic differential equations, 2022. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2207.06635)
670 [2207.06635](https://arxiv.org/abs/2207.06635).
- 671
- 672 Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation
673 using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference*
674 *on computer vision*, pp. 2223–2232, 2017a.
- 675
- 676 Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and
677 Eli Shechtman. Toward multimodal image-to-image translation. *Advances in neural information*
678 *processing systems*, 30, 2017b.
- 679 Shilong Zou, Yuhang Huang, Renjiao Yi, Chenyang Zhu, and Kai Xu. Cyclediff: Cycle diffusion
680 models for unpaired image-to-image translation. *arXiv preprint arXiv:2508.06625*, 2025.

682 7 APPENDIX

683 7.1 CONVERGENCE AND OPTIMALITY

684

685 To our knowledge, Zhu et al. (2017a) do not provide a proof of convergence for a cycle-consistent
686 framework. We provide a quick sketch below.

687

688 We begin by defining the spaces and mappings relevant to the cycle-consistent neural network frame-
689 work. Let X denote the input domain, where each element $x \in X$ represents one of n distinct items
690 (e.g., images), as determined by the dimensionality of the input data. Let Y be the target output
691 domain, where each $y \in Y$ likewise corresponds to one of n distinct items, consistent with the
692 output data dimension. Let $f : X \rightarrow Z$ denote the forward mapping implemented by the neural
693 network, where Z is an intermediate representation space. Let $g : Z \rightarrow X$ be the inverse mapping
694 used to reconstruct the original input from the network’s output. In the context of a model trained
695 with cycle-consistency loss, we assume the following two properties:

$$697 \quad \forall x \in X, \exists g \text{ such that } g(f(x)) = x \quad (10)$$

$$698 \quad Z = Y \quad (11)$$

699

700 Equation 10 ensures the existence of a cycle-consistent mapping. Equation 11 states that the inter-
701 mediate representation space Z is equivalent to the target output domain Y , thereby implying that
the network effectively learns a mapping from X to Y .

Given that f is bijective (by Equation 10) and that $|X| = |Y| = n$ (by Equation 11), it follows that there exist $n!$ possible one-to-one mappings (i.e., permutations) between elements of X and Y . Although many such bijective mappings exist in theory, the cycle-consistency loss biases the network toward converging on a single, consistent, and invertible transformation that minimizes the reconstruction error.

Let f_θ denote the forward neural network (parameterized by weights θ) which maps inputs from domain X to outputs in domain Y . Let g_θ denote the inverse neural network, also parameterized by θ , that attempts to reconstruct the original input from the output of f_θ . By Equation 10 and $|X| = |Y| = n$, note that $g = f^{-1}$. The cycle-consistency loss $\mathcal{L}(\theta)$ is then defined as the expected reconstruction error between the original input x and its reconstruction $g_\theta(f_\theta(x))$, measured using the squared L^2 norm:

$$\mathcal{L}(\theta) = \mathbb{E}_{x \sim X} \left[\|g_\theta(f_\theta(x)) - x\|_2^2 \right] \quad (12)$$

The optimal parameters θ^* are obtained by minimizing the cycle-consistency loss: $\theta^* = \arg \min_\theta \mathcal{L}(\theta)$.

In general, we assume the following about the solution (θ^*) landscape:

1. No local minima exist (i.e., network optimizer will never be stuck at a local minima)
2. There exists a unique θ^* such that $\mathcal{L}(\theta^*) < \epsilon$ for some $\epsilon \in \mathbb{R}^+$

Solution uniqueness is enforced by the neural network’s inherent incompleteness: the neural network cannot perfectly reconstruct x , i.e., $\mathcal{L}(\theta) > 0 \forall \theta$. Since exact recovery is impossible, the model cannot satisfy cycle-consistency for any parameterization/mapping. So, the model will choose the lowest θ^* for convergence. Under these assumptions, then, gradient descent will thus converge to a unique solution θ^* with corresponding invertible mappings ($f_{\theta^*}, g_{\theta^*}$) between domains X and Y .

Given a particular network architecture, then, if the human-preferred canonical map is suboptimal in perception-distortion, the model will converge to a different, possibly unusable solution. To avoid such behavior, small amounts of paired data can produce improved results.

7.2 TRAINING DETAILS

We provide a PyTorch implementation of our models.

The dataset consists of satellite photographs and images. We adopt train and test datasets from Zhu et al. (2017a), consisting of 1096 maps and satellite (aerial) photographs. Images are resized with random crop and flip to 256×256 , then normalized before training. For all probabilistic models, we set $\lambda_{cycle} = 10$, $\lambda_{id} = 5$, and $\lambda_{kl} = 1e - 05$ in Equation 9. We use the Adam optimizer (Kingma & Ba, 2015) with a batch size of 5. All networks were trained from scratch with a learning rate of 0.0002 and a latent dimension of $64 \times 16 \times 16$ (channels, height, and width respectively) for 600 epochs.

7.3 NETWORK ARCHITECTURE

Generator:

We use a U-Net style architecture (Johnson et al., 2016) with a variational bottleneck. The encoder and decoder are symmetric. Let $c7s1-k$ denote a 7×7 Convolution-InstanceNorm-ReLU layer with k filters and stride 1. dk denotes a custom Downsampling block via PixelUnshuffle with output channels k . Rk denotes a residual block that contains Reflection padding, two 3×3 convolutional layers with the same number of filters, two InstanceNorm layers and a ReLU activation.

Lk denotes a linear (1×1 convolutional) layer. Sk denotes a skip connection block that performs a linear projection or averaging to change the number of channels from the previous layer to k . uk denotes a custom Upsampling block via Pixelshuffle with output channels k . The network uses Reflection Padding throughout to reduce artifacts.

The model encodes an input image into parameters of a Gaussian distribution (mean, μ and log-variance, $\log \sigma$) in a learned latent space. The dimensionality of this space, N_z , determines the size of the bottleneck and the complexity of the learned representation. A latent code z is sampled

756 from this distribution using the reparameterization trick and is subsequently decoded to generate the
 757 output image.

758 For a 256×256 input image with 4 downsampling/upsampling layers and one residual block, the
 759 architecture is as follows:
 760

761 **Encoder:** $c7s1-64, d128, d256, d512, d1024, R1024, SN_z$
 762

763 **Variational Bottleneck:**

764 $\mu = S_{N_z}(S_{N_z}(\text{enc})), \log \sigma^2 = L_{N_z}(\text{enc}), \mathbf{z} \sim \mathcal{N}(\mu, \exp(\log \sigma^2))$
 765

766 **Decoder:** $S1024, R1024, u512, u256, u128, u64, c7s1-3$.

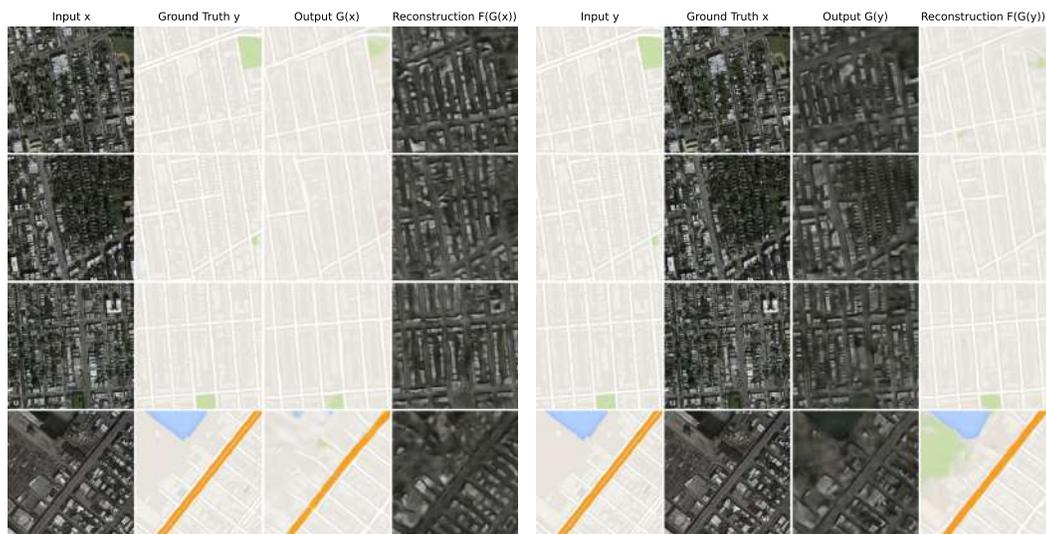
767 The final output layer ($c7s1-3$) consists of: ReflectionPad2d(3) \rightarrow Conv7-1 \rightarrow Tanh() activation.
 768

769 **Discriminator:**

770 The discriminator architecture is adopted from Zhu et al. (2017a), specifically utilizing a 70×70
 771 PatchGAN. The network is constructed from a series of 4×4 convolutional layers with Instance Nor-
 772 malization and LeakyReLU (slope 0.2). The first layer, C64 (64 filters, stride 2), omits Instance
 773 Normalization. This is followed by successive layers (C128, C256, C512), each doubling the num-
 774 ber of filters. A final convolution layer produces a 1-dimensional output map.
 775
 776
 777

778 7.4 TRANSLATION AND RECONSTRUCTION EXAMPLES

781 Each set of figures/images visualize (a) the aerial to map translation and aerial reconstruction and (b)
 782 map to aerial translation and map reconstruction. The maps translated by the Cycle-AE and Cycle-
 783 VAE models in an unpaired data setting show no structural similarity to the map domain, as there is
 784 neither an adversarial constraint nor a paired dataset to ensure domain alignment. However, spatial
 785 information is preserved and allows for the complete reconstruction of the original aerial input. We
 786 notice a similar pattern for aerial photos translation and map reconstruction.
 787
 788
 789

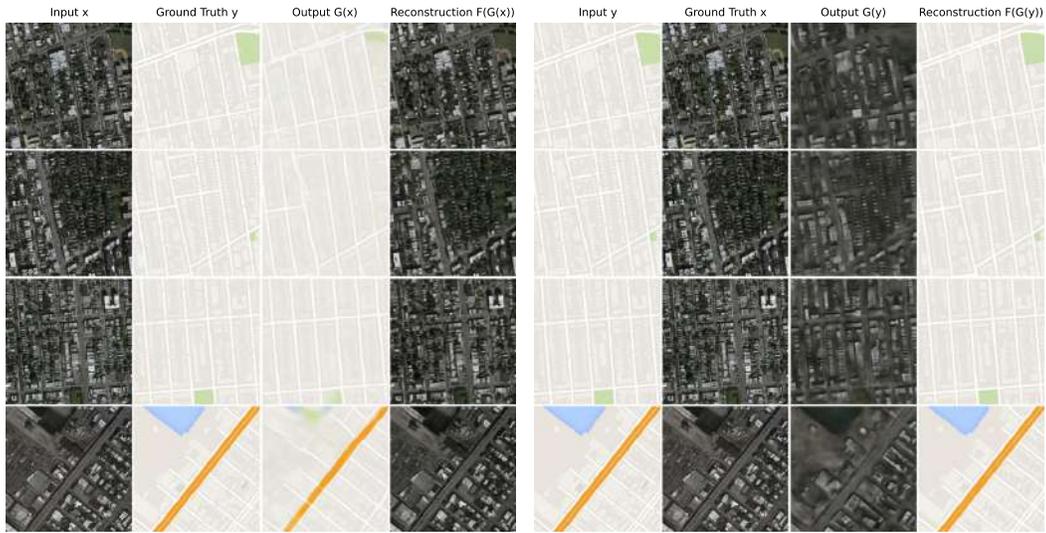


791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808 (a) Aerial to map translation and aerial reconstruction (b) Map to aerial translation and map reconstruction

809

Figure 12: AE translation and reconstruction.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863



(a) Aerial to map translation and aerial reconstruction (b) Map to aerial translation and map reconstruction

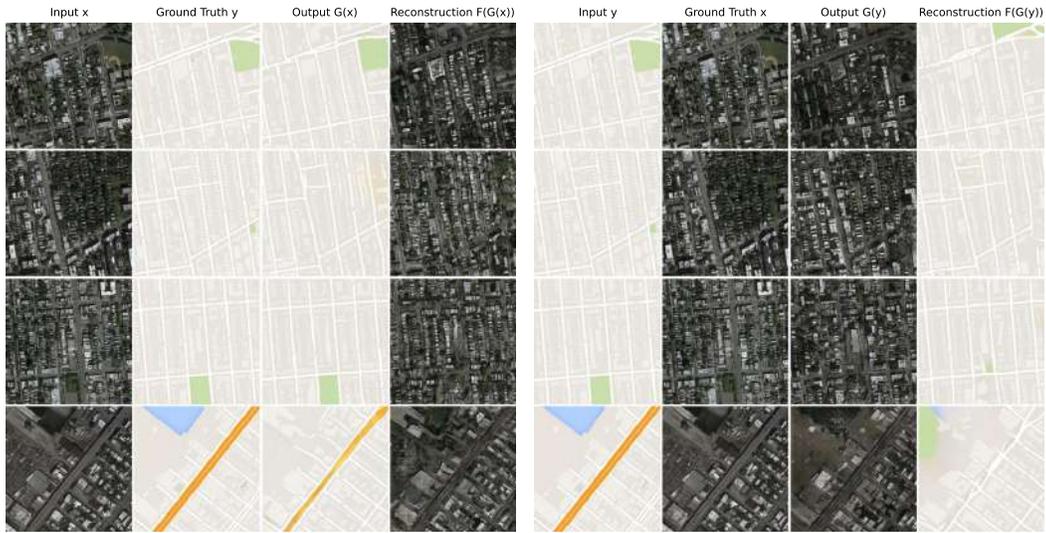
Figure 13: Cycle-AE (paired) translation and reconstruction.



(a) Aerial to map translation and aerial reconstruction (b) Map to aerial translation and map reconstruction

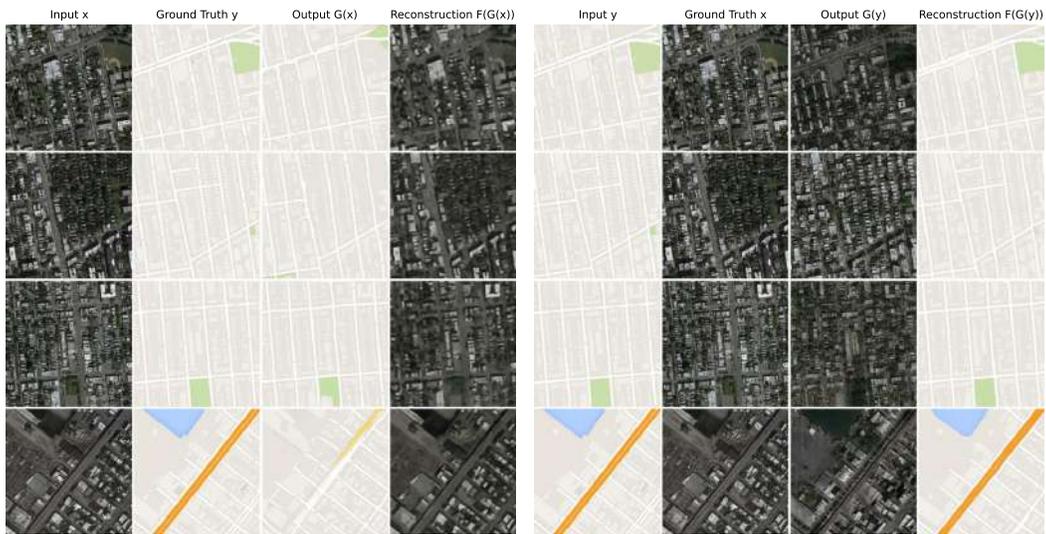
Figure 14: Cycle-AE (unpaired) translation and reconstruction.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917



(a) Aerial to map translation and aerial reconstruction (b) Map to aerial translation and map reconstruction

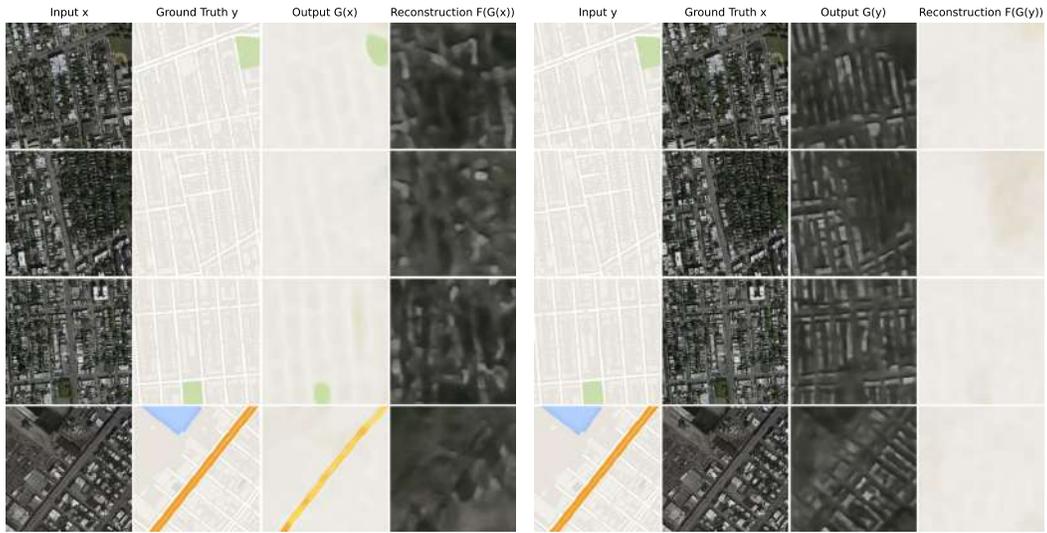
Figure 15: AE-GAN translation and reconstruction.



(a) Aerial to map translation and aerial reconstruction (b) Map to aerial translation and map reconstruction

Figure 16: AE CycleGAN translation and reconstruction.

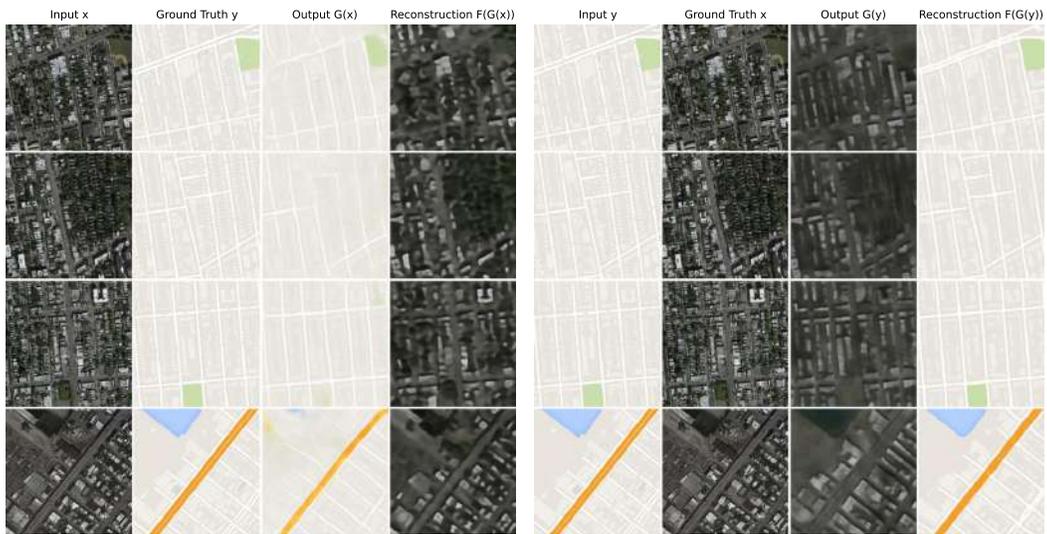
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936



(a) Aerial to map translation and aerial reconstruction (b) Map to aerial translation and map reconstruction

Figure 17: VAE translation and reconstruction.

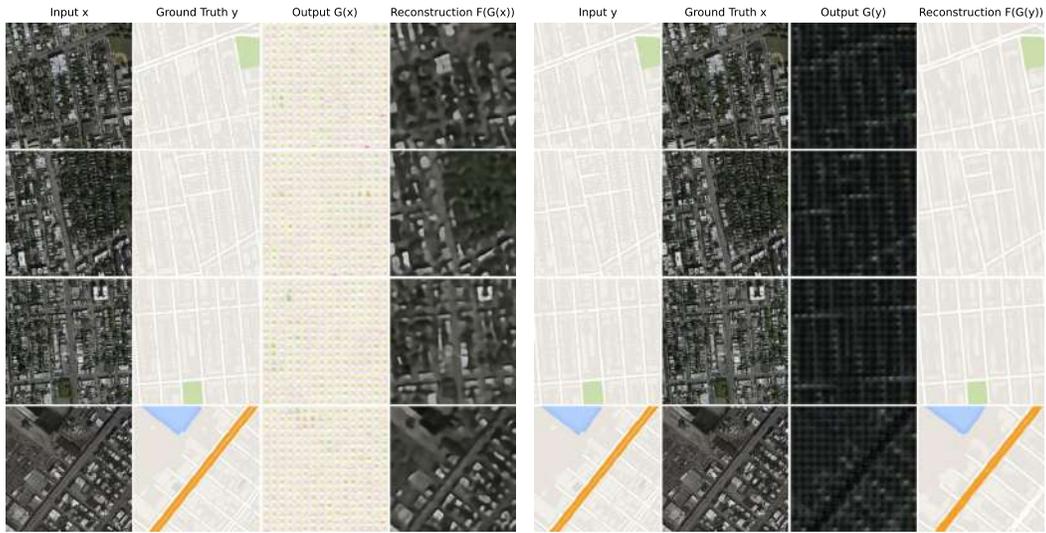
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971



(a) Aerial to map translation and aerial reconstruction (b) Map to aerial translation and map reconstruction

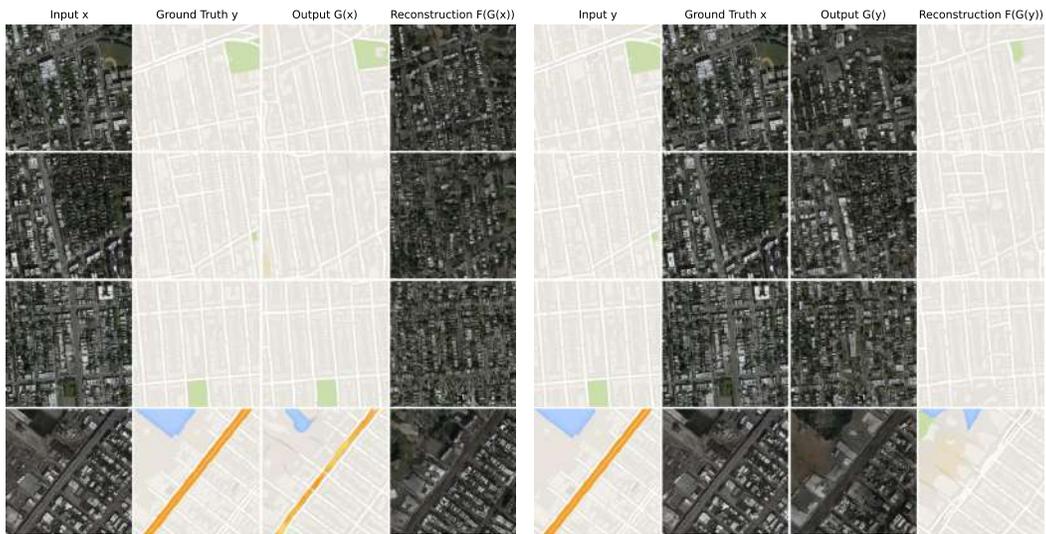
Figure 18: Cycle-VAE (paired) translation and reconstruction.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025



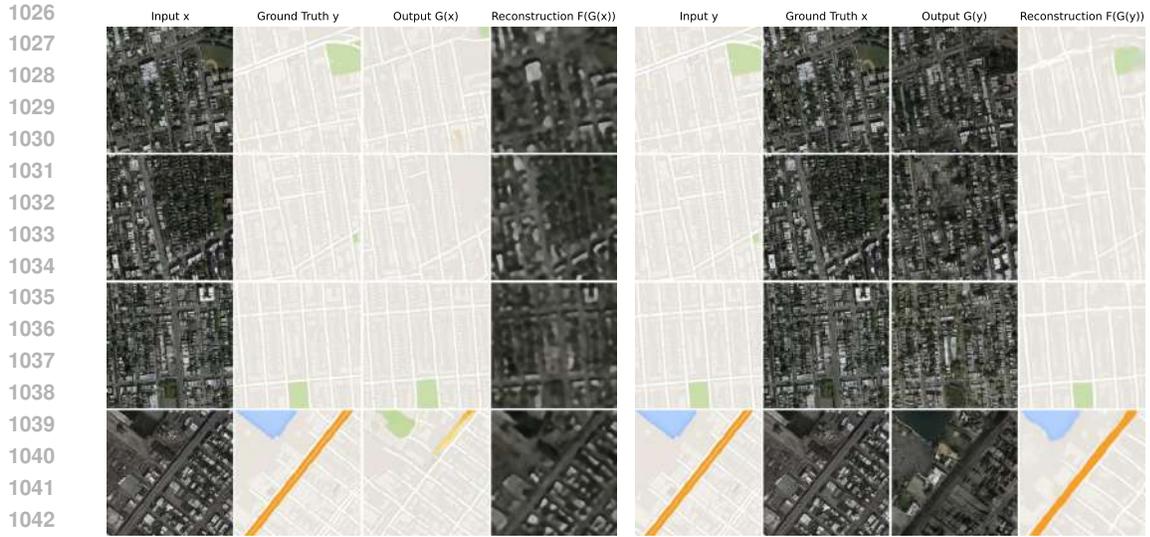
(a) Aerial to map translation and aerial reconstruction (b) Map to aerial translation and map reconstruction

Figure 19: Cycle-VAE (unpaired) translation and reconstruction.



(a) Aerial to map translation and aerial reconstruction (b) Map to aerial translation and map reconstruction

Figure 20: VAE-GAN translation and reconstruction.



(a) Aerial to map translation and aerial reconstruction (b) Map to aerial translation and map reconstruction

Figure 21: VAE CycleGAN translation and reconstruction.

7.4.1 ERROR EVALUATION METRICS

Table 5: Translation and reconstruction error evaluation of Autoencoder (AE) based models on aerial photos (x) \leftrightarrow maps (y) after 600 epochs

(a) AE						
AE	MSE	PSNR	SSIM	FID	KID Mean	KID STD
Aerial Translation	0.0256	14.08	0.2489	253.05	0.2331	0.0196
Map Translation	0.0018	27.37	0.7636	325.56	0.3530	0.0374
Aerial Reconstruction	0.0287	13.95	0.1992	288.35	0.2807	0.0208
Map Reconstruction	0.0017	27.77	0.7718	315.15	0.3347	0.0324
(b) CycleAE-paired						
CycleAE-paired	MSE	PSNR	SSIM	FID	KID Mean	KID STD
Aerial Translation	0.0251	14.13	0.2553	241.77	0.2175	0.0189
Map Translation	0.0018	27.46	0.7650	273.17	0.2729	0.0241
Aerial Reconstruction	0.0059	21.53	0.6661	175.29	0.1621	0.0164
Map Reconstruction	0.0003	35.91	0.9242	83.08	0.0604	0.0066
(c) CycleAE-unpaired						
CycleAE-unpaired	MSE	PSNR	SSIM	FID	KID Mean	KID STD
Aerial Translation	0.0410	13.69	0.0924	329.51	0.3552	0.0175
Map Translation	0.0872	10.59	0.0318	409.27	0.5077	0.0163
Aerial Reconstruction	0.0063	21.12	0.6451	180.12	0.1662	0.0188
Map Reconstruction	0.0002	36.59	0.9317	80.05	0.0608	0.0075
(d) AE-GAN						
AE-GAN	MSE	PSNR	SSIM	FID	KID Mean	KID STD
Aerial Translation	0.0349	14.29	0.2079	52.33	0.0175	0.0055
Map Translation	0.0031	25.10	0.6684	63.83	0.0316	0.0066
Aerial Reconstruction	0.0392	13.90	0.1481	70.75	0.0347	0.0067
Map Reconstruction	0.0050	22.99	0.6969	74.76	0.0425	0.0082
(e) AE-CycleGAN						
AE-CycleGAN	MSE	PSNR	SSIM	FID	KID Mean	KID STD
Aerial Translation	0.0389	13.62	0.1641	52.53	0.0170	0.0048
Map Translation	0.0050	22.97	0.6737	70.08	0.0460	0.0131
Aerial Reconstruction	0.0096	19.01	0.5069	200.11	0.1833	0.0181
Map Reconstruction	0.0005	32.99	0.8760	106.63	0.0844	0.0092

Table 6: Translation and reconstruction error evaluation of Variational Autoencoder (VAE) based models on aerial photos (x) \leftrightarrow maps (y) after 600 epochs.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

(a) VAE

VAE	MSE	PSNR	SSIM	FID	KID Mean	KID STD
Aerial Translation	0.0261	13.50	0.1894	304.34	0.2858	0.0209
Map Translation	0.0024	26.03	0.6753	358.70	0.3547	0.0310
Aerial Reconstruction	0.0340	11.13	0.1344	347.17	0.3642	0.0188
Map Reconstruction	0.0049	22.63	0.6717	367.13	0.3895	0.0250

(b) Cycle-VAE-paired

Cycle-VAE-paired	MSE	PSNR	SSIM	FID	KID Mean	KID STD
Aerial Translation	0.0255	13.25	0.2080	259.83	0.2403	0.0189
Map Translation	0.0022	26.58	0.7117	325.89	0.3369	0.0236
Aerial Reconstruction	0.0145	17.03	0.3330	266.76	0.2587	0.0173
Map Reconstruction	0.0007	31.81	0.8450	170.27	0.1546	0.0154

(c) Cycle-VAE-unpaired

Cycle-VAE-unpaired	MSE	PSNR	SSIM	FID	KID Mean	KID STD
Aerial Translation	0.0614	7.90	0.0885	475.98	0.5931	0.0274
Map Translation	0.0113	19.47	0.3128	419.51	0.5393	0.0194
Aerial Reconstruction	0.0146	17.31	0.3151	281.64	0.2761	0.0168
Map Reconstruction	0.0006	32.44	0.8468	219.37	0.2130	0.0282

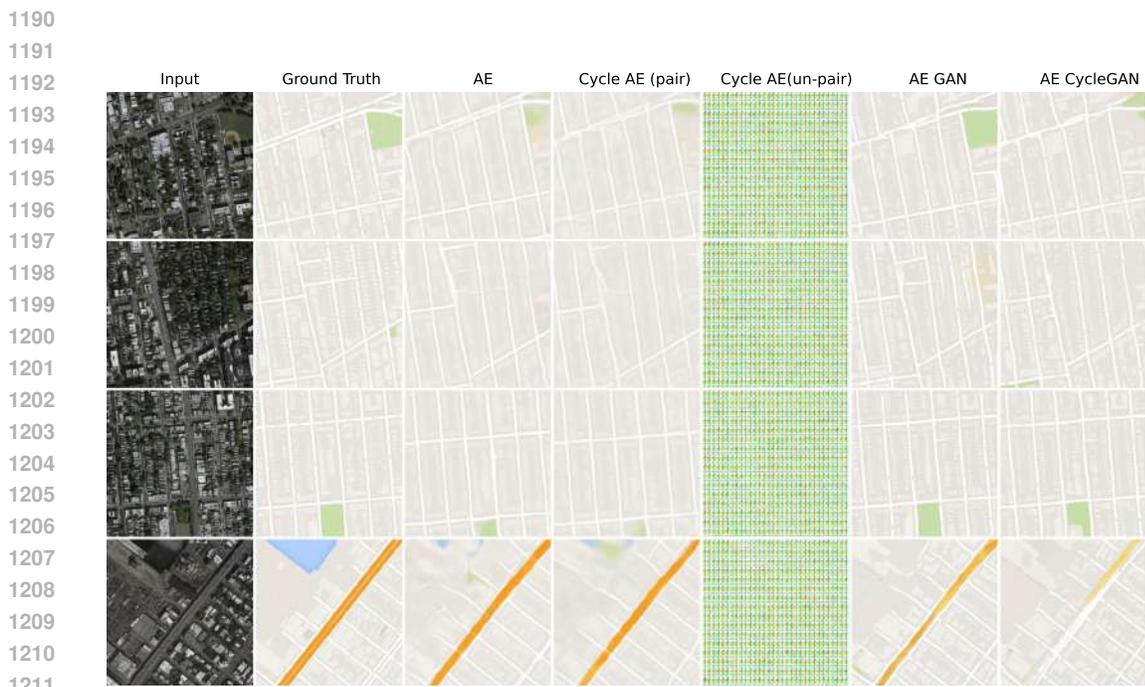
(d) VAE-GAN

VAE-GAN	MSE	PSNR	SSIM	FID	KID Mean	KID STD
Aerial Translation	0.0354	14.17	0.1705	64.45	0.0301	0.0052
Map Translation	0.0033	24.85	0.6272	83.18	0.0510	0.0091
Aerial Reconstruction	0.0430	13.35	0.0954	84.63	0.0496	0.0101
Map Reconstruction	0.0060	22.22	0.6172	104.49	0.0800	0.0107

(e) VAE-CycleGAN

VAE-CycleGAN	MSE	PSNR	SSIM	FID	KID Mean	KID STD
Aerial Translation	0.0443	13.36	0.0965	69.25	0.0378	0.0063
Map Translation	0.0056	22.50	0.5793	90.87	0.0594	0.0092
Aerial Reconstruction	0.0175	16.10	0.2779	271.72	0.2703	0.0200
Map Reconstruction	0.0011	29.67	0.7804	241.78	0.2409	0.0259

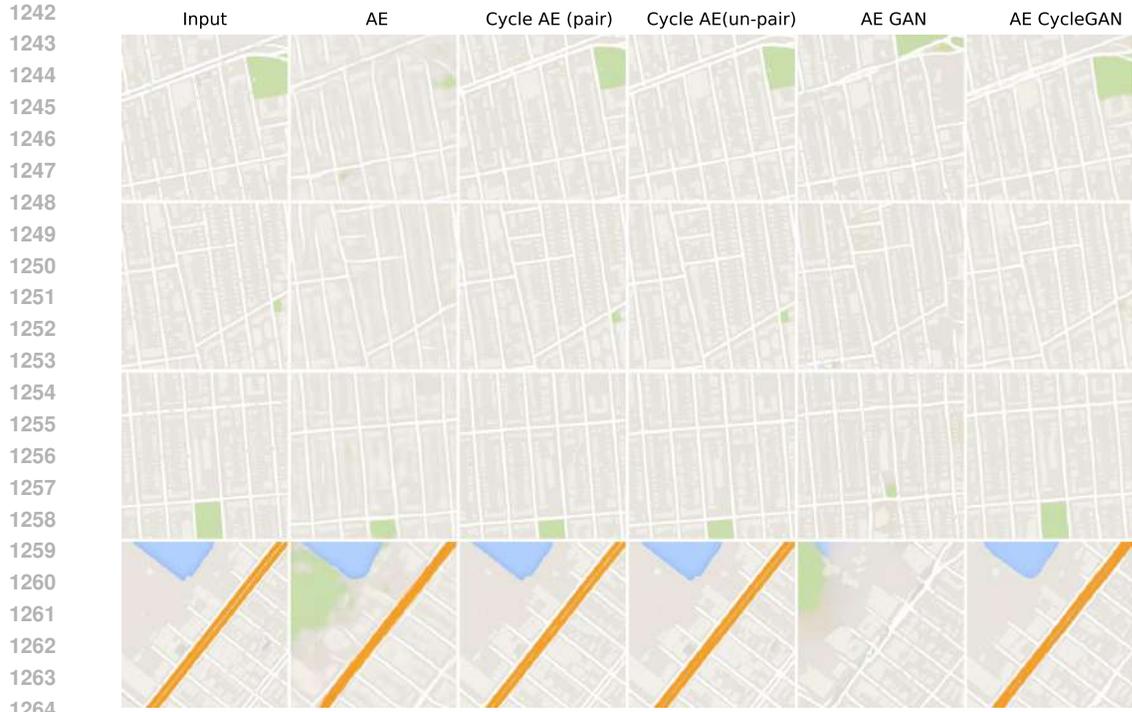
1188 7.4.2 COMPARISON: AE VARIANTS
 1189
 1190
 1191



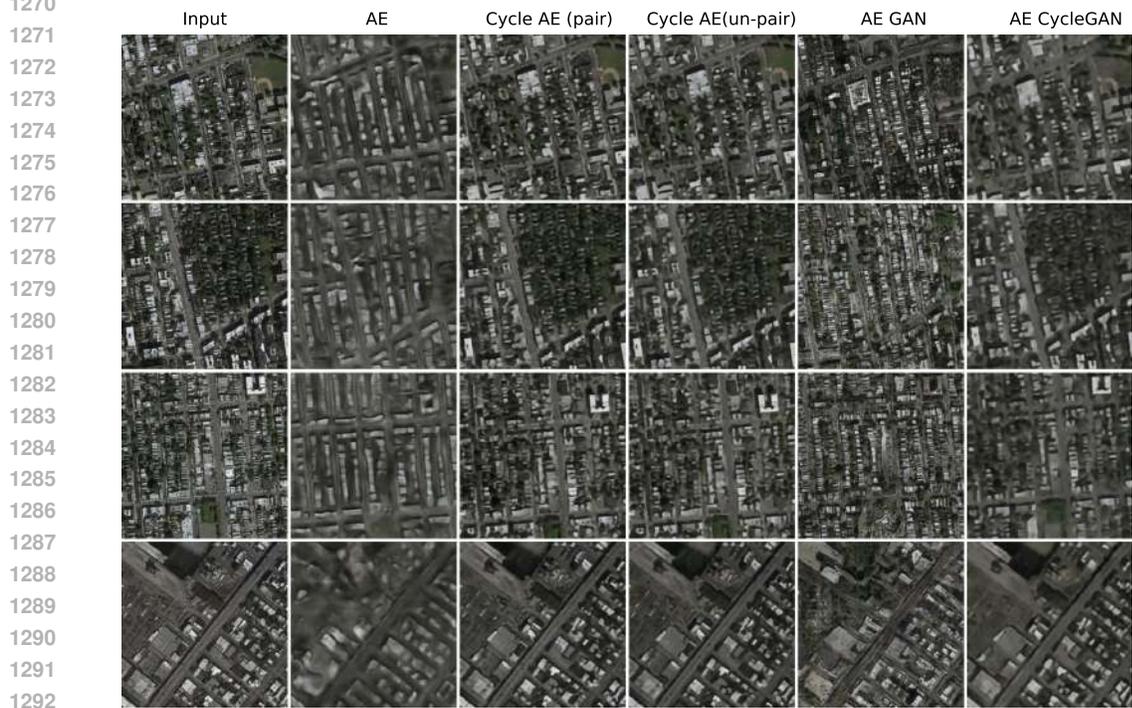
1212 Figure 22: Translated maps $x \rightarrow G(x)$ of different AE models for the given aerial input x . From
 1213 left to right: input, ground truth, translated map output from models AE, Cycle-AE (paired data),
 1214 Cycle-AE (unpaired data), AE-GAN, and AE-CycleGAN.
 1215



1239 Figure 23: Translated aerial images $y \rightarrow F(y)$ of different AE models for the given map input y .
 1240 From left to right: input, ground truth, translated aerial images from models AE, Cycle-AE (paired
 1241 data), Cycle-AE (unpaired data), AE-GAN, and AE-CycleGAN.

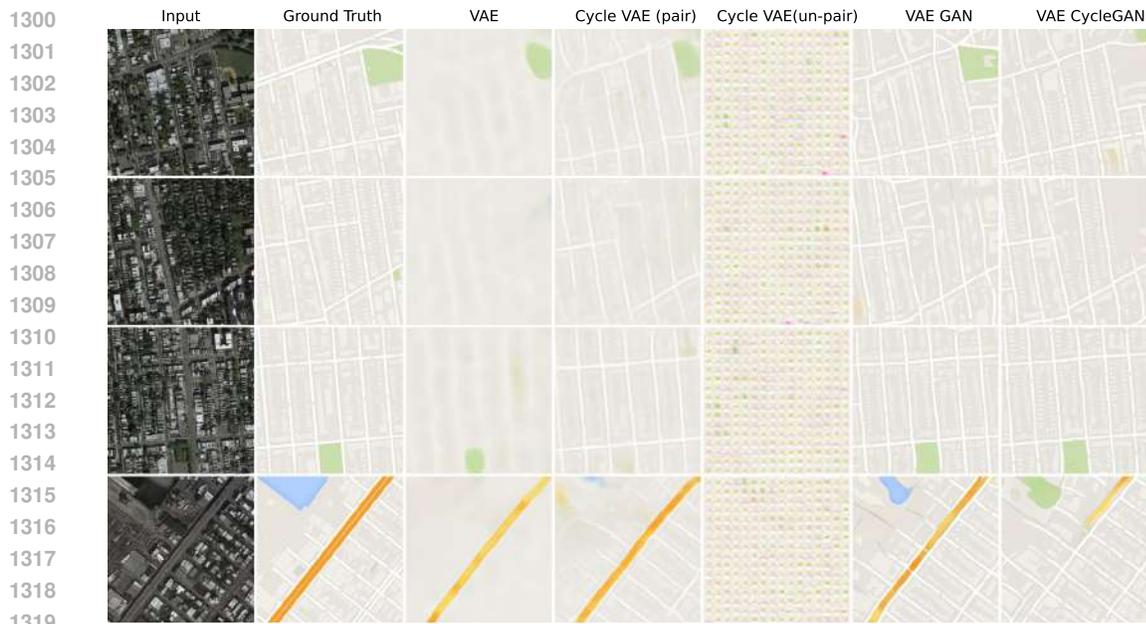


1265 Figure 24: Reconstructed maps $G(F(y))$ of different AE models for the given input y . From left to
1266 right: input y , reconstructed maps from the models AE, Cycle-AE (paired data), Cycle-AE (unpaired
1267 data), AE-GAN, and AE-CycleGAN.
1268



1293 Figure 25: Reconstructed aerial images $F(G(x))$ of different AE models for the given input x . From
1294 left to right: input x , reconstructed maps from the models AE, Cycle-AE (paired data), Cycle-AE
1295 (unpaired data), AE-GAN, and AE-CycleGAN.

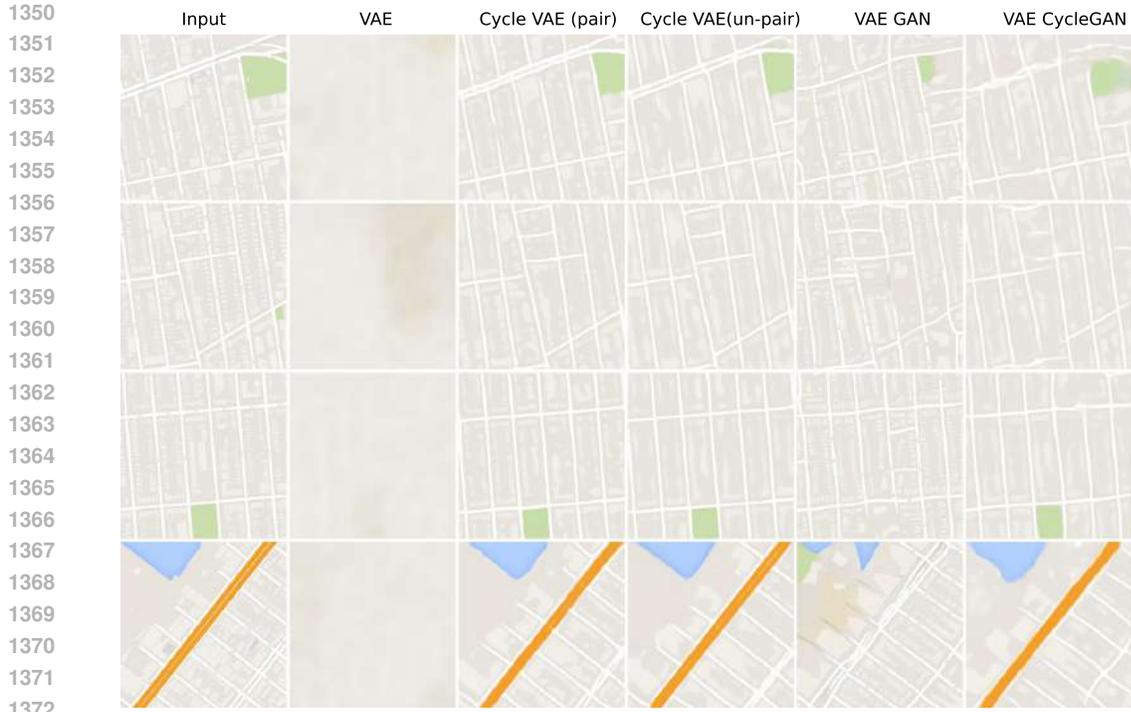
1296 7.4.3 COMPARISON: VAE VARIANTS
 1297
 1298
 1299



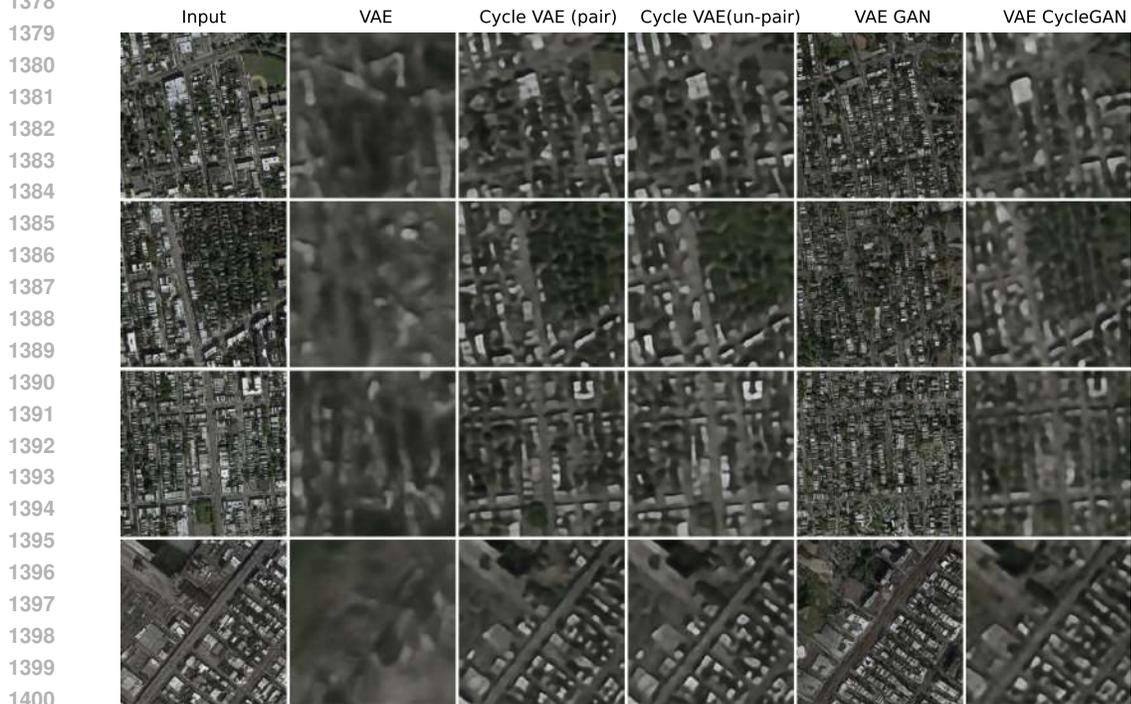
1320 Figure 26: Translated maps $x \rightarrow G(x)$ of different VAE models for the given aerial input x . From
 1321 left to right: input, ground truth, translated map output from models VAE, Cycle-VAE (paired data),
 1322 Cycle-VAE (unpaired data), VAE-GAN, and VAE-CycleGAN.
 1323
 1324
 1325
 1326



1347 Figure 27: Translated aerial images $y \rightarrow F(y)$ of different VAE models for the given map input
 1348 y . From left to right: input, ground truth, translated aerial images from models VAE, Cycle-VAE
 1349 (paired data), Cycle-VAE (unpaired data), VAE-GAN, and VAE-CycleGAN.

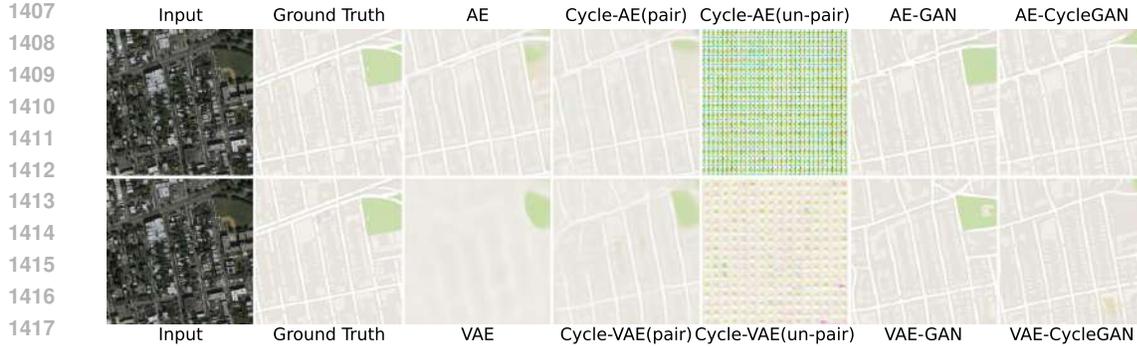


1373 Figure 28: Reconstructed maps $G(F(y))$ of different VAE models for the given input y . From left
 1374 to right: input y , reconstructed maps from the models VAE, Cycle-VAE (paired data), Cycle-VAE
 1375 (unpaired data), VAE-GAN, and VAE-CycleGAN.
 1376

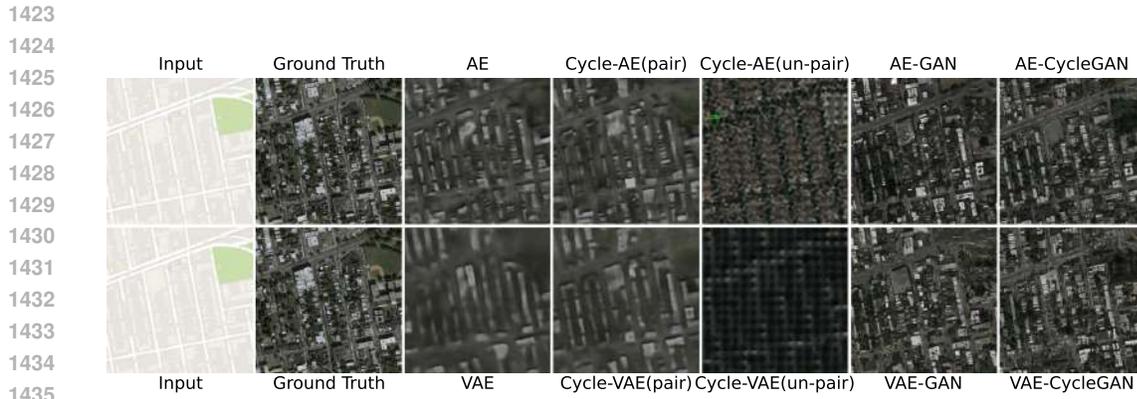


1401 Figure 29: Reconstructed aerial images $F(G(x))$ of different VAE models for the given input x .
 1402 From left to right: input x , reconstructed aerial images from the models VAE, Cycle-VAE (paired
 1403 data), Cycle-VAE (unpaired data), VAE-GAN, and VAE-CycleGAN.

1404 7.4.4 COMPARISON BETWEEN AE AND VAE MODELS
 1405
 1406



1419 Figure 30: Comparison of translated map outputs $G(x)$ from different models for a given input x .
 1420 From left to right: input x , ground truth y , translated map output from the models. From top to
 1421 bottom: AE variants, VAE variants.



1436 Figure 31: Comparison of translated aerial image outputs $F(y)$ from different models for a given
 1437 input y . From left to right: input y , ground truth x , translated aerial output from the models. From
 1438 top to bottom: AE variants, VAE variants.



1456 Figure 32: Comparison of reconstructed maps $G(F(y))$ from different models for a given input y .
 1457 From left to right: input y , reconstructed map output from the models. From top to bottom: AE
 variants, VAE variants.

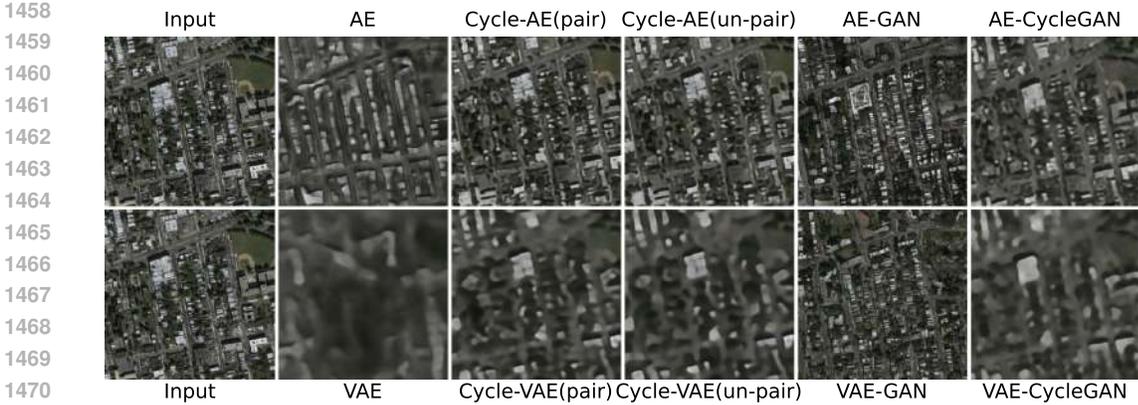


Figure 33: Comparison of reconstructed aerial image outputs $F(G(x))$ from different models for a given input y . From left to right: input x , reconstructed aerial output from the models. From top to bottom: AE variants, VAE variants.

7.4.5 ERROR EVALUATION: TRANSLATION AND RECONSTRUCTION

Table 7: Translation error evaluation: comparison between deterministic and stochastic models.

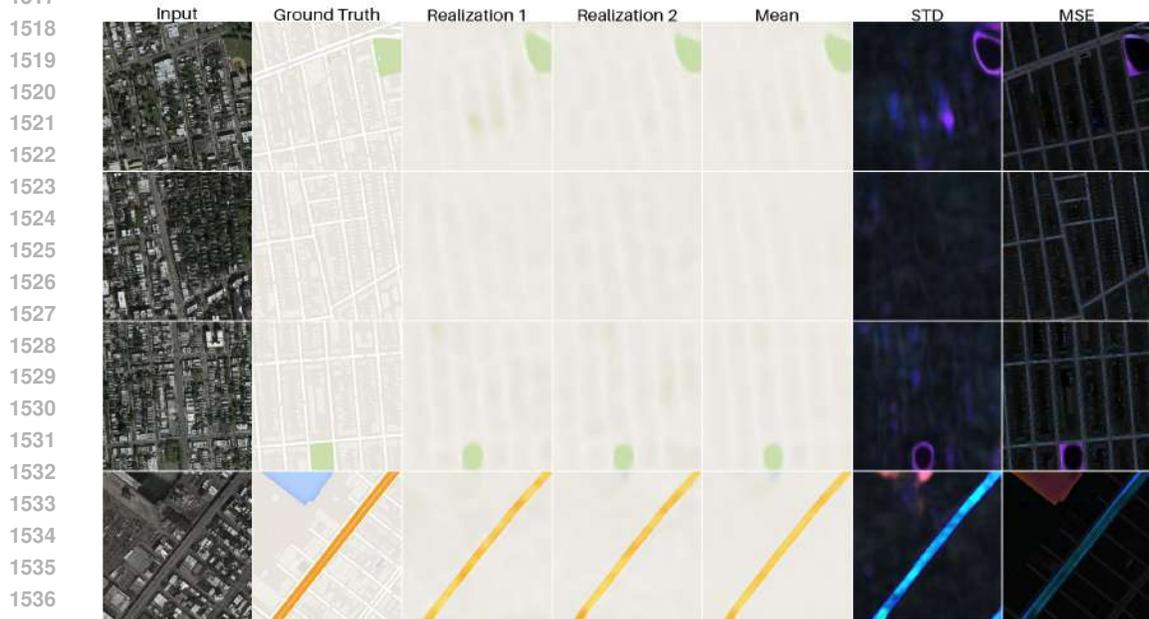
Model	$G : X \rightarrow Y$ (aerial \rightarrow map)					$F : Y \rightarrow X$ (map \rightarrow aerial)					
	MSE \downarrow	PSNR \uparrow	SSIM \uparrow	FID \downarrow	KID $\mu \pm \sigma \downarrow$	MSE \downarrow	PSNR \uparrow	SSIM \uparrow	FID \downarrow	KID $\mu \pm \sigma \downarrow$	
Deterministic	AE	0.0018	27.37	0.7636	325.56	0.3530 \pm 0.0374	0.0256	14.08	0.2489	253.05	0.2331 \pm 0.0196
	Cycle AE (paired)	0.0018	27.46	0.7650	273.17	0.2729 \pm 0.0241	0.0251	14.13	0.2553	241.77	0.2175 \pm 0.0189
	Cycle AE (unpaired)	0.0872	10.59	0.0318	409.27	0.5077 \pm 0.0163	0.0410	13.69	0.0924	329.51	0.3552 \pm 0.0175
	AE-GAN	0.0031	25.10	0.6684	63.83	0.0316 \pm 0.0066	0.0349	14.29	0.2079	52.33	0.0175 \pm 0.0055
AE-CycleGAN	0.0050	22.97	0.6737	70.08	0.0460 \pm 0.0131	0.0389	13.62	0.1641	52.53	0.0170 \pm 0.0048	
Stochastic	VAE	0.0024	26.03	0.6753	358.70	0.3547 \pm 0.0310	0.0261	13.50	0.1894	304.34	0.2858 \pm 0.0209
	Cycle VAE (paired)	0.0022	26.58	0.7117	325.89	0.3369 \pm 0.0236	0.0255	13.25	0.2080	259.83	0.2403 \pm 0.0189
	Cycle VAE (unpaired)	0.0113	19.47	0.3128	419.51	0.5393 \pm 0.0194	0.0614	7.90	0.0885	475.98	0.5931 \pm 0.0274
	VAE-GAN	0.0033	24.85	0.6272	83.18	0.0510 \pm 0.0091	0.0354	14.17	0.1705	64.45	0.0301 \pm 0.0052
VAE-CycleGAN	0.0056	22.50	0.5793	90.87	0.0594 \pm 0.0092	0.0443	13.36	0.0965	69.25	0.0378 \pm 0.0063	

Table 8: Reconstruction error: comparison between deterministic and stochastic models.

Model	$G(F(y))$: Map Reconstruction					$F(G(x))$: Aerial Reconstruction					
	MSE \downarrow	PSNR \uparrow	SSIM \uparrow	FID \downarrow	KID $\mu \pm \sigma \downarrow$	MSE \downarrow	PSNR \uparrow	SSIM \uparrow	FID \downarrow	KID $\mu \pm \sigma \downarrow$	
Deterministic	AE	0.0017	27.77	0.7718	315.15	0.3347 \pm 0.0324	0.0287	13.95	0.1992	288.35	0.2807 \pm 0.0208
	Cycle AE (paired)	0.0003	35.91	0.9242	83.08	0.0604 \pm 0.0066	0.0059	21.53	0.6661	175.29	0.1621 \pm 0.0164
	Cycle AE (unpaired)	0.0002	36.59	0.9317	80.05	0.0608 \pm 0.0075	0.0063	21.12	0.6451	180.12	0.1662 \pm 0.0188
	AE-GAN	0.0050	22.99	0.6969	74.76	0.0425 \pm 0.0082	0.0392	13.90	0.1481	70.75	0.0347 \pm 0.0067
AE-CycleGAN	0.0005	32.99	0.8760	106.63	0.0844 \pm 0.0092	0.0096	19.01	0.5069	200.11	0.1833 \pm 0.0181	
Stochastic	VAE	0.0049	22.63	0.6717	367.13	0.3895 \pm 0.0250	0.0340	11.13	0.1344	347.17	0.3642 \pm 0.0188
	Cycle VAE (paired)	0.0007	31.81	0.8450	170.27	0.1546 \pm 0.0154	0.0145	17.03	0.3330	266.76	0.2587 \pm 0.0173
	Cycle VAE (unpaired)	0.0006	32.44	0.8468	219.37	0.2130 \pm 0.0282	0.0146	17.31	0.3151	281.64	0.2761 \pm 0.0168
	VAE-GAN	0.0060	22.22	0.6172	104.49	0.0800 \pm 0.0107	0.0430	13.35	0.0954	84.63	0.0496 \pm 0.0101
VAE-CycleGAN	0.0011	29.67	0.7804	241.78	0.2409 \pm 0.0259	0.0175	16.10	0.2779	271.72	0.2703 \pm 0.0200	

1512 7.5 REALIZATIONS
1513

1514 7.5.1 VARIATIONAL AUTOENCODER (VAE) REALIZATIONS
1515
1516



1538 Figure 34: VAE model map realizations. From left to right: input (x), ground truth (y), 2 sample
1539 realizations, mean, STD, MSE of the 30 map realizations. STD is scaled by $20\times$ for visibility.
1540

1541

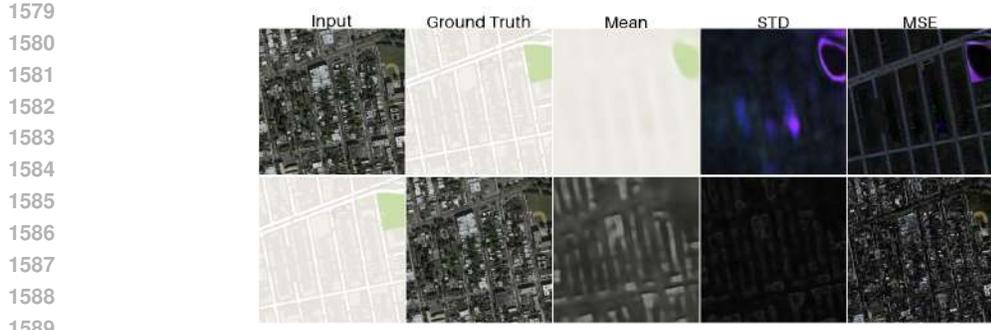
1542



1563 Figure 35: VAE model aerial photo realizations. From left to right: input (y), ground truth (x), 2
1564 sample realizations, mean, STD and MSE of the 30 aerial realizations. STD is scaled by $10\times$ for
1565 visibility.

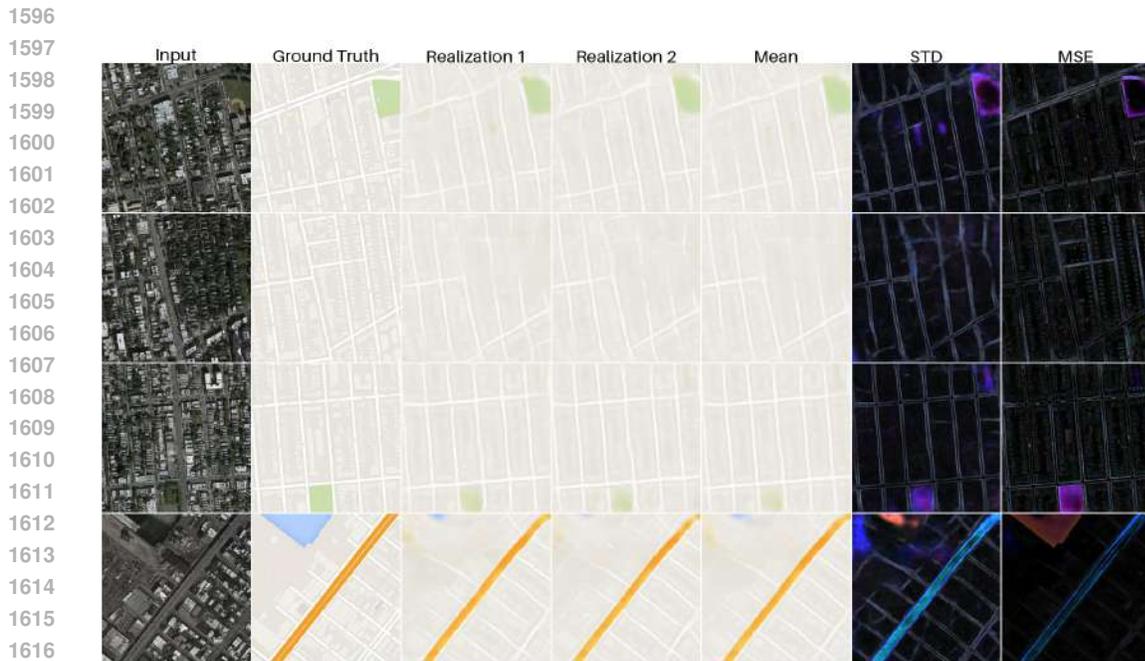


1577 Figure 36: Cycle VAE (paired) sample aerial and map realizations.



1590 Figure 37: Cycle-VAE (paired) model mean, STD and MSE of the 30 realizations. Top row: aerial \rightarrow
1591 maps, bottom row: maps \rightarrow aerial images. For visibility, the STD of the aerial image is scaled by
1592 10x and the STD of the map image by 20x.

1593
1594
1595 **7.5.2 CYCLE-VAE (PAIRED) REALIZATIONS**

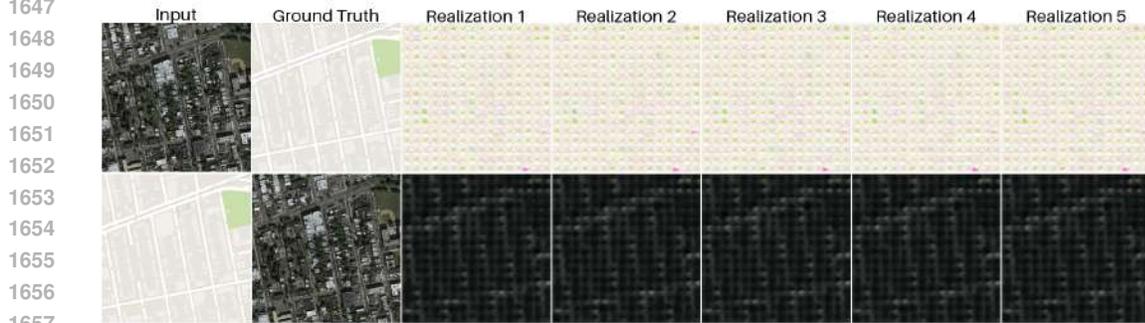


1618 Figure 38: Cycle-VAE (paired) model map realizations. From left to right: input (x), ground truth
1619 (y), 2 sample realizations, mean, STD, MSE of the 30 map realizations. STD is scaled by 20x for
visibility.

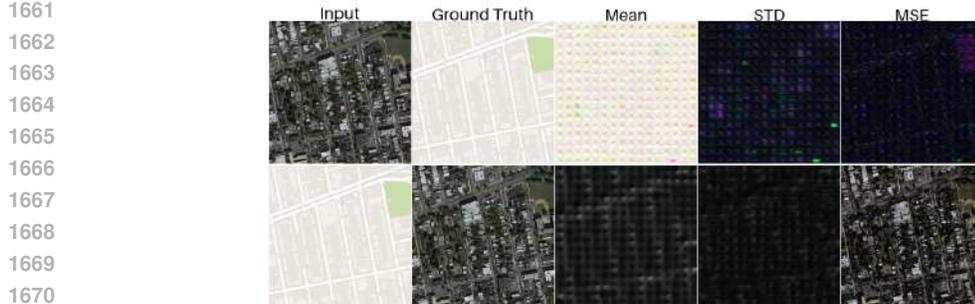


1640
1641
1642
1643
1644
Figure 39: Cycle-VAE (paired) model aerial photo realizations. From left to right: input (y), ground truth (x), 2 sample realizations, mean, STD and MSE of the 30 aerial realizations. STD is scaled by $10\times$ for visibility.

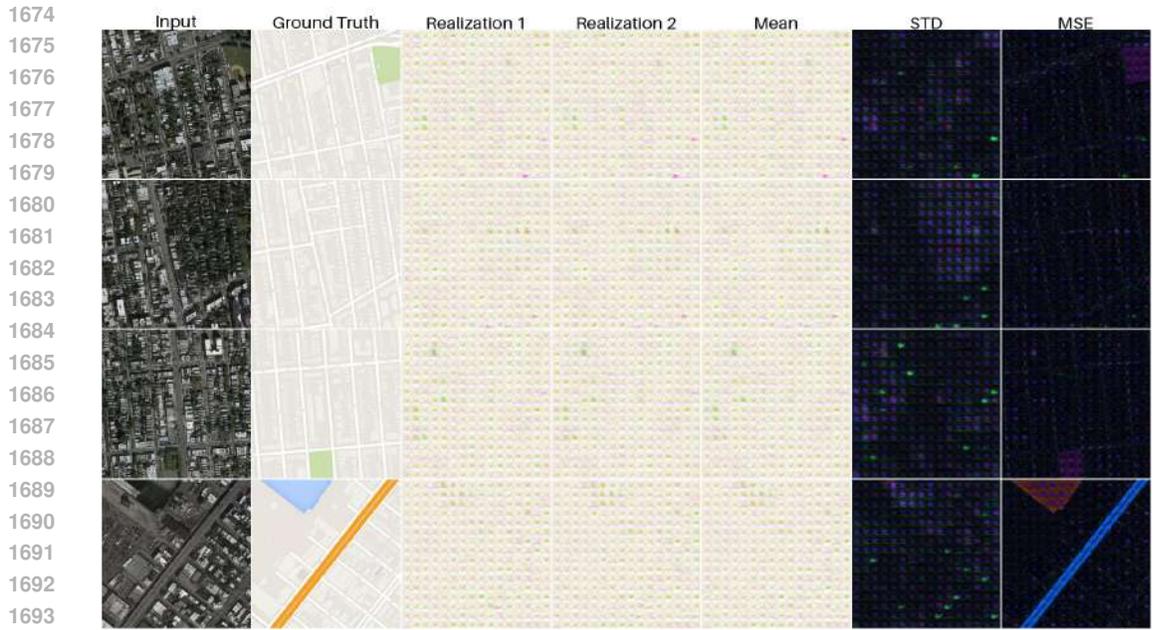
1645
1646
7.5.3 CYCLE-VAE (UNPAIRED) REALIZATIONS



1658
1659
1660
Figure 40: Cycle VAE (unpaired) sample aerial and map realizations.



1671
1672
1673
Figure 41: Cycle-VAE (unpaired) model mean, STD and MSE of the 30 realizations. Top row: aerial \rightarrow maps, bottom row: maps \rightarrow aerial images. For visibility, the STD of the aerial image is scaled by $10\times$ and the STD of the map image by $20\times$.



1695
1696
1697
1698
1699
1700
1701
1702

Figure 42: Cycle-VAE (unpaired) model map realizations. From left to right: input (x), ground truth (y), 2 sample realizations, mean, STD, MSE of the 30 map realizations. STD is scaled by $20\times$ for visibility.



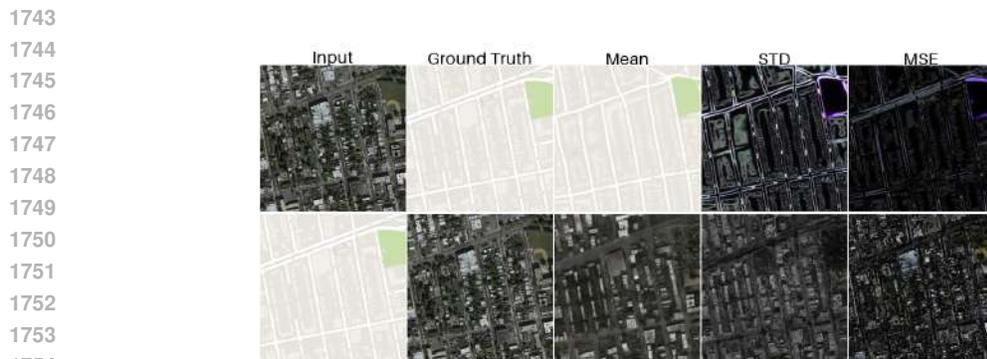
1723
1724
1725
1726
1727

Figure 43: Cycle-VAE (unpaired) model aerial photo realizations. From left to right: input (y), ground truth (x), 2 sample realizations, mean, STD, and MSE of the 30 aerial realizations. STD is scaled by $10\times$ for visibility.

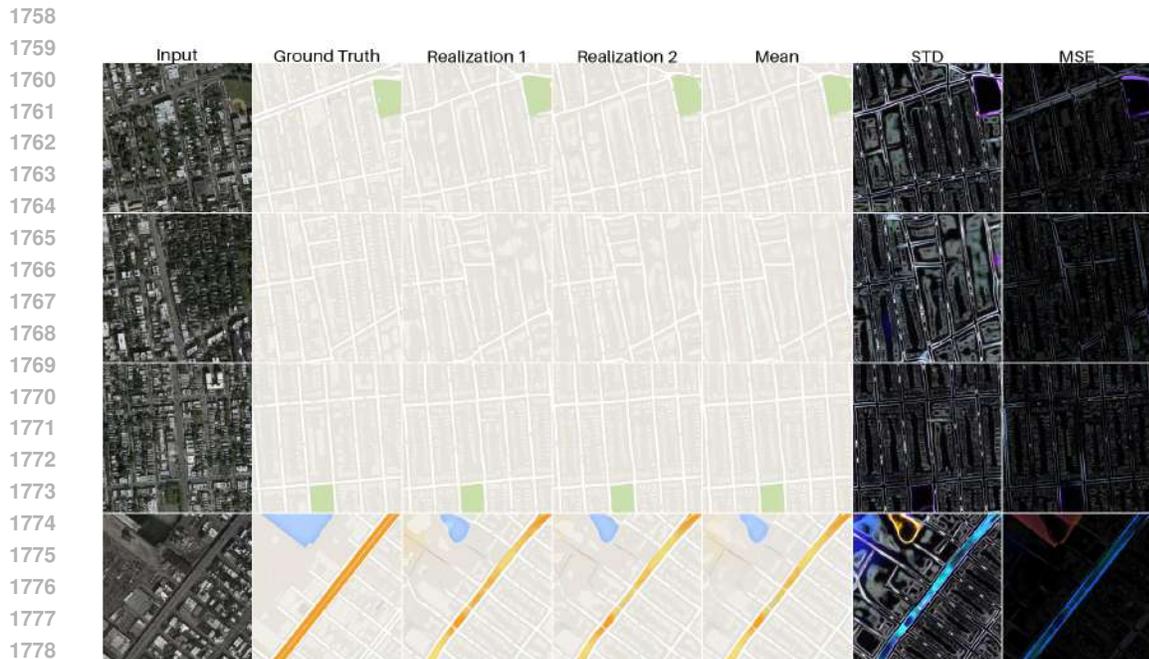
1728 7.5.4 VAE-GAN REALIZATIONS
 1729



1741 Figure 44: VAE-GAN sample aerial and map realizations.
 1742



1755 Figure 45: VAE-GAN model mean, STD and MSE of the 30 realizations. Top row: aerial \rightarrow maps,
 1756 bottom row: maps \rightarrow aerial images. For visibility, the STD of the aerial image is scaled by $3\times$ and
 1757 the STD of the map image by $20\times$.

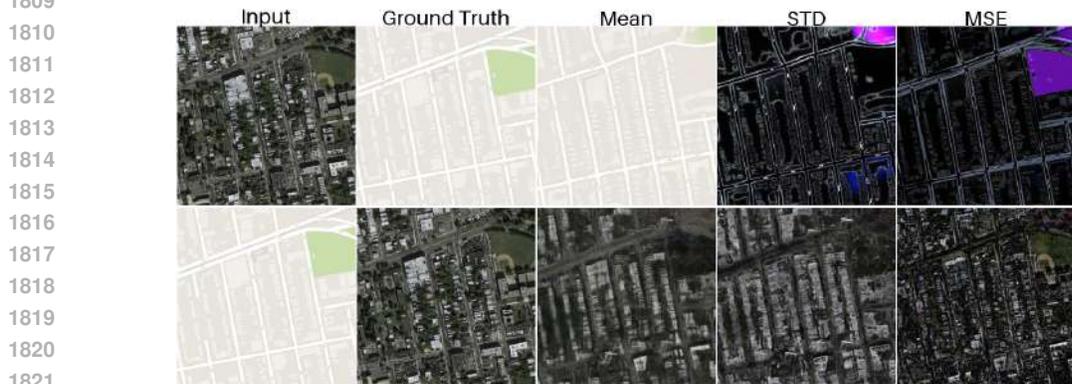


1780 Figure 46: VAE-GAN model aerial photo realizations. From left to right: input (y), ground truth
 1781 (x), 2 sample realizations, mean, STD and MSE of the 30 aerial realizations. STD is scaled by $3\times$
 for visibility.



1803 Figure 47: VAE-GAN model aerial photo realizations. From left to right: input (y), ground truth
1804 (x), 2 sample realizations, mean, STD and MSE of the 30 aerial realizations. STD is scaled by $3\times$
1805 for visibility.

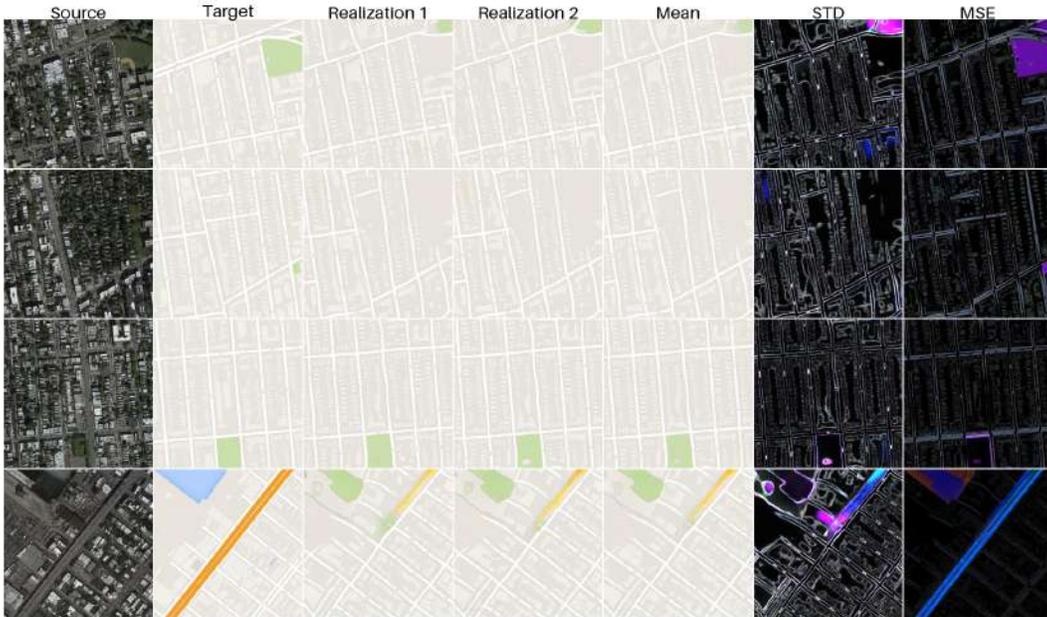
1806
1807 7.5.5 VAE-CYCLEGAN REALIZATIONS



1822 Figure 48: Mean, scaled STD (aerial $3\times$, maps $20\times$), and MSE across 30 realizations.

1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857



1858 Figure 49: VAE-CycleGAN model map realizations. From left to right: input (x), ground truth (y),
1859 2 sample realizations, mean, STD and MSE of the 30 map realizations. STD is scaled by $25\times$ for
1860 visibility.
1861
1862
1863
1864

1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884



1885
1886 Figure 50: VAE-CycleGAN model aerial photo realizations. From left to right: input (y), ground
1887 truth (x), 2 sample realizations, mean, STD and MSE of the 30 aerial realizations. STD is scaled by
1888 $3\times$ for visibility.
1889

7.6 REALIZATION METRICS

Table 9: MSE comparisons for 30 realizations, evaluated on aerial photos (x) \leftrightarrow maps (y). As the true translation posterior is intractable, high MSE may either indicate mean error or a (correct) region of high variance in the translation.

Table 10: VAE

VAE	MSE (y, \hat{y}_i) Maps	MSE (x, \hat{x}_i) Aerial photos
Image 1	0.001517	0.024690
Image 2	0.001211	0.023881
Image 3	0.001258	0.027295
Image 4	0.006711	0.023370
Average (30)	0.006711	0.023370

Table 11: Cycle VAE (paired)

Cycle VAE (paired)	MSE (y, \hat{y}_i) Maps	MSE (x, \hat{x}_i) Aerial photos
Image 1	0.001212	0.023442
Image 2	0.000959	0.025219
Image 3	0.001035	0.027436
Image 4	0.006002	0.022164
Average (30)	0.002302	0.024565

Table 12: Cycle VAE (unpaired)

Cycle VAE (unpaired)	MSE (y, \hat{y}_i) Maps	MSE (x, \hat{x}_i) Aerial photos
Image 1	0.008954	0.061927
Image 2	0.007704	0.056705
Image 3	0.008576	0.071990
Image 4	0.023110	0.053937
Average (30)	0.012086	0.061140

Table 13: VAE-GAN

VAE-GAN	MSE (y, \hat{y}_i) Maps	MSE (x, \hat{x}_i) Aerial photos
Image 1	0.001687	0.028285
Image 2	0.001156	0.030473
Image 3	0.000963	0.031132
Image 4	0.009378	0.025543
Average (30)	0.003296	0.028858

Table 14: VAE-CycleGAN

VAE-CycleGAN	MSE (y, \hat{y}_i) Maps	MSE (x, \hat{x}_i) Aerial photos
Image 1	0.003122	0.033576
Image 2	0.001444	0.036811
Image 3	0.001688	0.036307
Image 4	0.019514	0.033555
Average (30)	0.006442	0.035062

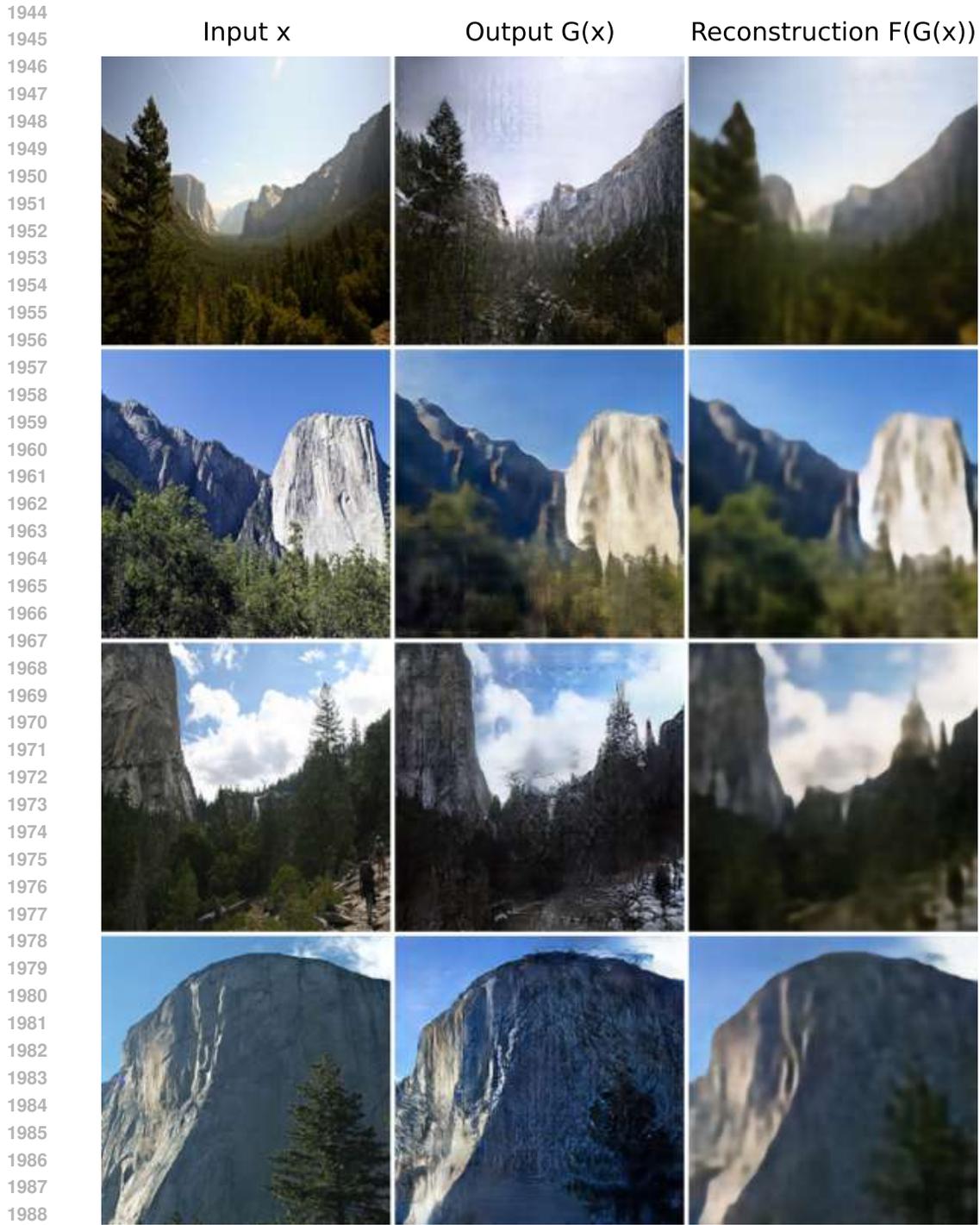


Figure 51: Summer \leftrightarrow Winter : VAE-CycleGAN translated and reconstructed images, $X \rightarrow G(X) \rightarrow F(G(X))$

1990
1991
1992
1993
1994
1995
1996
1997

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

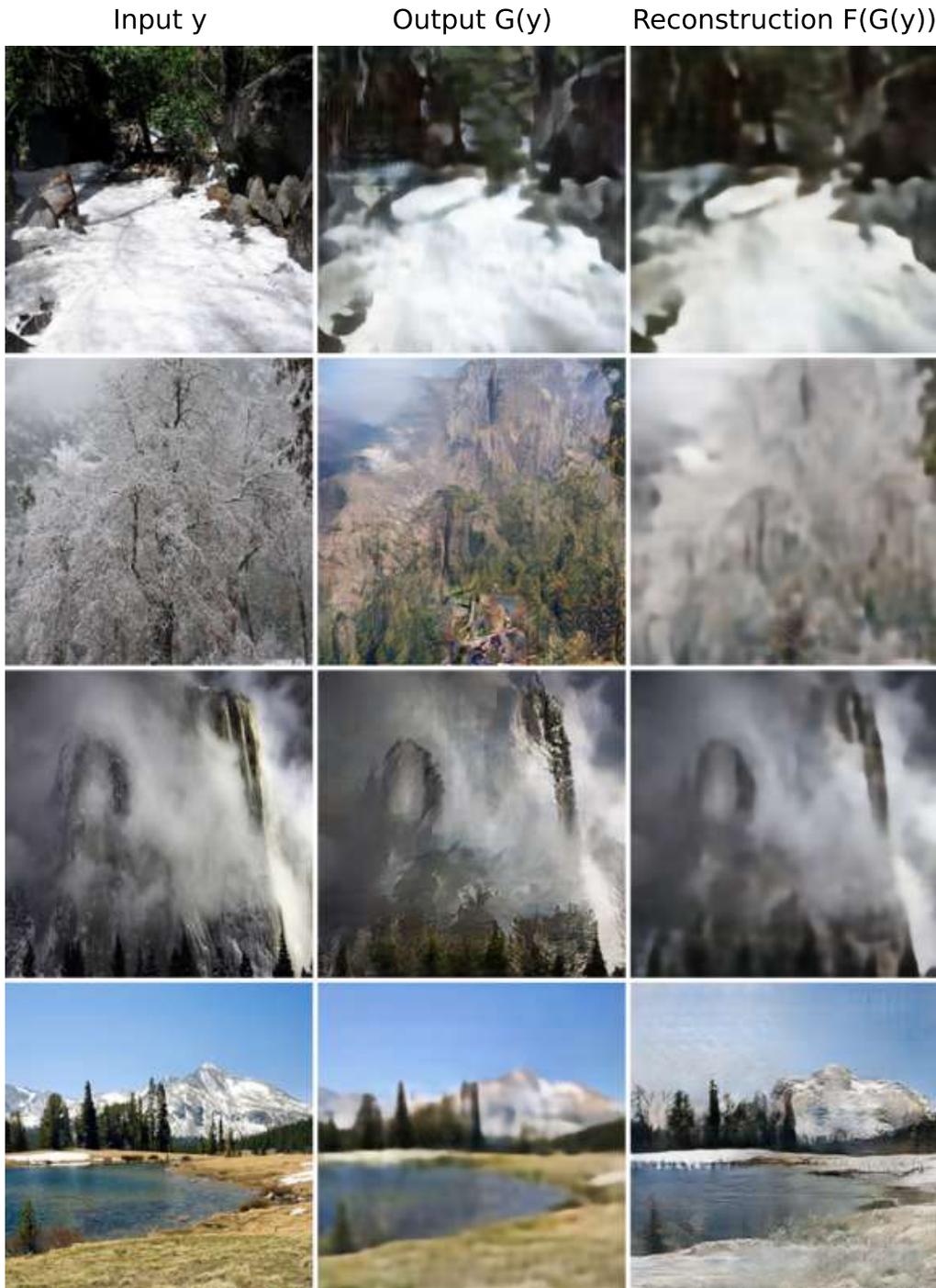


Figure 52: Summer \leftrightarrow Winter : VAE-CycleGAN translated and reconstructed images, $Y \rightarrow F(Y) \rightarrow G(F(Y))$

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

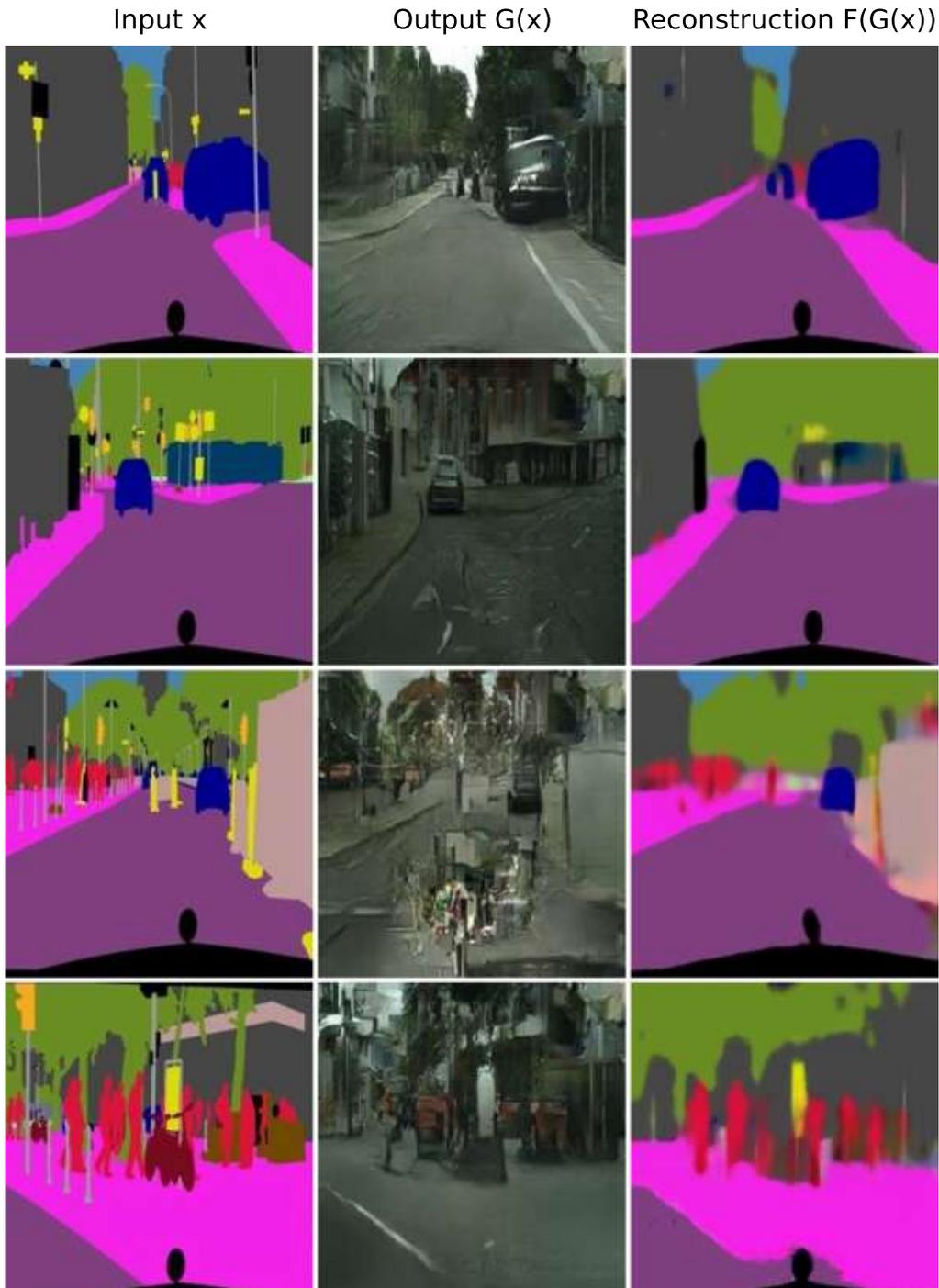


Figure 53: Label \leftrightarrow Cityscape: VAE-CycleGAN translated and reconstructed images, $X \rightarrow G(X) \rightarrow F(G(X))$

2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159



Figure 54: Label \leftrightarrow Cityscape: VAE-CycleGAN translated and reconstructed images, $Y \rightarrow F(Y) \rightarrow G(F(Y))$