

WHO SAID NEURAL NETWORKS AREN'T LINEAR?

Nimrod Berman*

Ben-Gurion University

bermann@post.bgu.ac.il

Assaf Hallak*

NVIDIA

ahallak@nvidia.com

Assaf Shocher*†

Technion

assaf.sh@technion.ac.il

ABSTRACT

Neural networks are famously nonlinear. However, linearity is defined relative to a pair of vector spaces, $f: \mathcal{X} \rightarrow \mathcal{Y}$. Is it possible to identify a pair of non-standard vector spaces for which a conventionally nonlinear function is, in fact, linear? This paper introduces a method that makes such vector spaces explicit by construction. We find that if we sandwich a linear operator A between two invertible neural networks, $f(x) = g_y^{-1}(Ag_x(x))$, then the corresponding vector spaces \mathcal{X} and \mathcal{Y} are induced by newly defined addition and scaling actions derived from g_x and g_y . We term this kind of architecture a Linearizer. This framework makes the entire arsenal of linear algebra, including SVD, pseudo-inverse, orthogonal projection and more, applicable to nonlinear mappings. Furthermore, we show that the composition of two Linearizers that share a neural network is also a Linearizer. We leverage this property and demonstrate that training diffusion models using our architecture makes the hundreds of sampling steps collapse into a single step. We further utilize our framework to enforce idempotency (i.e. $f(f(x)) = f(x)$) on networks leading to a globally projective generative model and to demonstrate modular style transfer.[‡]

1 INTRODUCTION

Linearity occupies a privileged position in mathematics, physics, and engineering. Linear systems admit a rich and elegant theory: they can be decomposed through eigenvalue and singular value analysis, inverted or pseudo-inverted with well-understood stability guarantees, and manipulated compositionally without loss of structure. These properties are not only aesthetically pleasing, they underpin the computational efficiency of countless algorithms in signal processing, control theory, and scientific computing. Crucially, repeated application of linear operators simplifies rather than complicates: iteration reduces to powers of eigenvalues, continuous evolution is captured by the exponential of an operator, and composition preserves linearity (Strang, 2022).

In contrast, non-linear systems, while more expressive, often defy such a structure. Iterating non-linear mappings can quickly lead to intractable dynamics; inversion may be ill-posed or undefined; and even simple compositional questions lack closed-form answers (Strogatz, 2024). Neural networks, the dominant modeling tool in modern machine learning, are famously nonlinear, placing their analysis and manipulation outside the reach of classical linear algebra (Hornik et al., 1989). As a result, tasks that are trivial in the linear setting, such as projecting onto a subspace, enforcing idempotency, or collapsing iterative procedures, become major challenges in the nonlinear regime that require engineered loss functions and optimization schemes.

This raises a natural and fundamental question: can we reinterpret conventionally nonlinear neural networks as linear operators in disguise? If so, we would gain access to the full arsenal of linear methods without sacrificing the expressive power of nonlinear function classes.

In this paper, we propose a framework that precisely achieves this goal. By embedding linear operators between two invertible neural networks, we induce new vector space structures under which the overall mapping is linear. We call such architectures *Linearizers*. This perspective not only offers a new lens on neural networks, but also enables powerful applications: collapsing hundreds of diffusion sampling steps into one, enforcing structural properties such as idempotency, and more. In short, Linearizers provide a bridge between the expressive flexibility of nonlinear models and the analytical tractability of linear algebra.

*Equal contribution.

†A.S. is a Chaya Fellow, supported by the Chaya Career Advancement Chair.

[‡]Code available at <https://github.com/assafshocher/Linearizer>.

2 LINEARIZER FRAMEWORK

Linearity is not an absolute concept; it is a property defined relative to a pair of vector spaces, $f : \mathcal{X} \rightarrow \mathcal{Y}$. This begs the question: is it possible to identify a pair of spaces \mathcal{X} and \mathcal{Y} , for which a function that is nonlinear over standard Euclidean space is, in fact, perfectly linear? To answer this, we must first deconstruct what defines a vector space. It is a structure comprised of a set of vectors (e.g. \mathbb{R}^N), a field of scalars (e.g., \mathbb{R}), and two fundamental operations: vector addition (+) and scalar multiplication (\cdot). We propose that while keeping the set of vectors and the field unchanged, we can redefine the operations in order to induce new vector spaces.

2.1 DEFINITIONS

We introduce a formalism we term as *Linearizer*, in which the relevant vector spaces are immediately identifiable by construction. The architecture we propose can be trained from scratch or by distillation from an existing model. Our approach is made practical by invertible neural networks Rezende & Mohamed (2015); Dinh et al. (2015; 2017). We build our model by wrapping a linear operator, a matrix A , between two such invertible networks, g_x and g_y :

Definition 1 (Linearizer). *Let \mathcal{X}, \mathcal{Y} be two spaces and $g_x : \mathcal{X} \rightarrow \mathcal{X}$, $g_y : \mathcal{Y} \rightarrow \mathcal{Y}$ be two corresponding invertible functions. Also, let $A : \mathcal{X} \rightarrow \mathcal{Y}$ be a linear operator. Then we define the Linearizer $\mathbb{L}_{\{g_x, g_y, A\}}$ as the following function $f : \mathcal{X} \rightarrow \mathcal{Y}$:*

$$f(x) = \mathbb{L}_{\{g_x, g_y, A\}}(x) = g_y^{-1}(Ag_x(x)) \quad (1)$$

For this construction, we can define the corresponding vector spaces \mathcal{X} and \mathcal{Y} , also shown in Figure 1:

Definition 2 (Induced Vector Space Operations). *Let $g : V \rightarrow V$ be an invertible function. We define a new set of operations, \oplus and \odot , for any two vectors $v_1, v_2 \in V$ and any scalar $a \in \mathbb{R}$:*

$$v_1 \oplus_g v_2 := g^{-1}(g(v_1) + g(v_2)) \quad (2)$$

$$a \odot_g v_1 := g^{-1}(a \cdot g(v_1)) \quad (3)$$

2.2 LINEARITY

The input space \mathcal{X} is defined by operations (\oplus_x, \odot_x) induced by g_x , and the output space \mathcal{Y} by (\oplus_y, \odot_y) induced by g_y .

Lemma 1. (V, \oplus_g, \odot_g) is a vector space over \mathbb{R} . *Proof.* See Appendix A.

Lemma 2. *The function $f(x)$ is a linear map from the vector space \mathcal{X} to the vector space \mathcal{Y} .*

Proof. To prove linearity, we show that f satisfies the superposition principle.

$$\begin{aligned} f(a_1 \odot_x x_1 \oplus_x a_2 \odot_x x_2) &= g_y^{-1}(Ag_x(g_x^{-1}(a_1 g_x(x_1) + a_2 g_x(x_2)))) \\ &= g_y^{-1}(a_1 g_y(g_y^{-1}(Ag_x(x_1))) + a_2 g_y(g_y^{-1}(Ag_x(x_2)))) \\ &= a_1 \odot_y f(x_1) \oplus_y a_2 \odot_y f(x_2) \end{aligned} \quad \square$$

2.3 INTUITION

The Linearizer can be understood through several intuitive lenses.

Linear algebra analogy. The Linearizer is analogous to an eigendecomposition or SVD. Just as a matrix is diagonal (i.e. simplest) in its eigenbasis, our function f is a simple matrix multiplication in the "linear basis" defined by g_x and g_y . This provides a natural coordinate system in which to analyze and manipulate the function. For example, repeated application of a Linearizer with a shared basis ($g_x = g_y = g$) is equivalent to taking a power of its core matrix:

$$\mathbb{L}_{\{g, g, A\}}^{\circ N}(x) = \underbrace{f(f(\dots f(x)))}_{N \text{ times}} = g^{-1}(A^N \cdot g(x)) = \mathbb{L}_{\{g, g, A^N\}}(x) \quad (4)$$

Geometric interpretation. g_x can be seen as a diffeomorphism. Equip the input with the pullback metric $g_x^* \langle \cdot, \cdot \rangle_{\mathbb{R}^N}$ induced by the Euclidean metric in latent space. With respect to this metric, straight lines in latent correspond to geodesics in input, and linear interpolation in the induced coordinates maps to smooth and semantically coherent curves in input space.

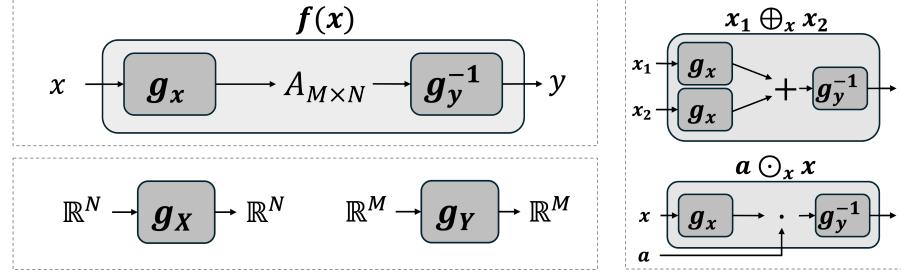


Figure 1: **Left.** The Linearizer structure (top) is a linear operation sandwiched between two invertible functions (bottom) **Right.** Vector addition and scalar multiplication define induced vector spaces for which f is linear.

2.4 PROPERTIES

We review how basic properties of linear transforms are expressed by Linearizers.

2.4.1 Composition

Lemma 3. *The composition of two Linearizer functions with compatible spaces $f_1 : \mathcal{X} \rightarrow \mathcal{Y}$ and $f_2 : \mathcal{Y} \rightarrow \mathcal{Z}$, is also a Linearizer.*

$$\text{Proof. } (f_2 \circ f_1)(x) = g_z^{-1}(A_2 g_y(g_y^{-1}(A_1 g_x(x)))) = g_z^{-1}((A_2 A_1) \cdot g_x(x)) \quad \square \quad (5)$$

2.4.2 Inner Product and Hilbert Space

Definition 3 (Induced Inner Product). *Given an invertible map $g : V \rightarrow \mathbb{R}^N$, the induced inner product on V is*

$$\langle v_1, v_2 \rangle_g := \langle g(v_1), g(v_2) \rangle_{\mathbb{R}^N}, \quad (6)$$

where the right-hand side is the standard Euclidean dot product.

Equipped with this inner product, the induced vector spaces make *Hilbert spaces*.

Lemma 4 (Induced spaces are Hilbert). *Let $g : V \rightarrow \mathbb{R}^n$ be a bijection and endow V with the induced vector space operations (\oplus_g, \odot_g) and inner product*

$$\langle u, v \rangle_g := \langle g(u), g(v) \rangle_{\mathbb{R}^n}. \quad (7)$$

Then $(V, \langle \cdot, \cdot \rangle_g)$ is a Hilbert space. Proof. See Appendix C

2.4.3 Transpose

Lemma 5 (Transpose). *Let $f(x) = g_y^{-1}(Ag_x(x))$ be a Linearizer. Its transpose $f^\top : \mathcal{Y} \rightarrow \mathcal{X}$ with respect to the induced inner products is*

$$f^\top(y) = g_x^{-1}(A^\top g_y(y)). \quad (8)$$

Proof. For all $x \in \mathcal{X}, y \in \mathcal{Y}$,

$$\langle f(x), y \rangle_{g_y} = \langle Ag_x(x), g_y(y) \rangle_{\mathbb{R}^N} = \langle g_x(x), A^\top g_y(y) \rangle_{\mathbb{R}^N} \quad (9)$$

$$= \langle x, g_x^{-1}(A^\top g_y(y)) \rangle_{g_x} = \langle x, f^\top(y) \rangle_{g_x}. \quad \square \quad (10)$$

2.4.4 Singular Value Decomposition

Lemma 6 (SVD of a Linearizer). *Let $A = U\Sigma V^\top$ be the singular value decomposition of A . Then the SVD of $f(x) = g_y^{-1}(Ag_x(x))$ is given by singular values Σ , input singular vectors $\tilde{v}_i = g_x^{-1}(v_i)$, and output singular vectors $\tilde{u}_i = g_y^{-1}(u_i)$.*

Proof. See Appendix B. \square

3 APPLICATIONS

Having established the theoretical foundations, we now present proof-of-concept applications that highlight advantages of our framework, leaving scaling and engineering refinements for future work.

3.1 ONE-STEP FLOW MATCHING

Flow matching (Lipman et al., 2023; Song et al., 2021b) trains a network to predict the velocity field that transports noise samples to data points. Traditionally, this velocity must be integrated over time with many small steps, making generation slow. We train flow matching using our Linearizer as the back bone, with a single vector space ($g_x = g_y = g$). The key property that helps us achieve one-step generation is the closure of linear operators. Applying the model sequentially with some simple linear operations in between steps is a chain of linear operations. This chain can be collapsed to a single linear operator so that the entire trajectory is realized in a single step.

Training. Training is just as standard FM, only using the Linearizer model and the induced spaces. We define the forward diffusion process:

$$x_t = (1-t) \odot x_0 \oplus t \odot x_1 = g^{-1}((1-t)g(x_0) + t g(x_1)). \quad (11)$$

This is a straight line in the g -space, mapped back to a curve in the data space. The target velocity is

$$v = x_1 \ominus x_0 = g^{-1}(g(x_1) - g(x_0)). \quad (12)$$

We parameterize a time-dependent Linearizer

$$f(x_t, t) = g^{-1}(A_t g(x_t)), \quad (13)$$

and train it to predict v . The loss is

$$\mathcal{L} = \mathbb{E}_{x_0, x_1, t} \| v \ominus f(x_t, t) \|^2 = \mathbb{E}_{x_0, x_1, t} \| g^{-1}(g(x_1) - g(x_0) - A_t g(x_t)) \|^2. \quad (14)$$

This objective ensures that the learned operators A_t approximate the true velocity of the flow within the latent space of the Linearizer.

Sampling as a collapsed operator. In inference time, standard practice is to discretize the ODE with many steps. In the induced space, an Euler update reads

$$x_{t+\Delta t} = x_t \oplus (\Delta t \odot f(x_t, t)), \quad (15)$$

which expands to

$$g(x_{t+\Delta t}) = (I + \Delta t A_t) g(x_t). \quad (16)$$

Iterating this N times produces

$$g(\hat{x}_1) \approx \underbrace{\left[\prod_{i=0}^{N-1} (I + \Delta t A_{t_i}) \right]}_{:=B} g(x_0). \quad (17)$$

The entire product can be collected into a single operator B , resulting in

$$\hat{x}_1 = g^{-1}(B g(x_0)). \quad (18)$$

Thus, what was originally a long sequence of updates collapses into a single multiplication. In practice, higher-order integration such as Runge-Kutta (Runge, 1895; Kutta, 1901) may be used to calculate B more accurately (see Appendix F). B is calculated only once, after training. Then generation requires only a single feedforward activation of our model on noise, regardless of the number of steps discretizing the ODE.

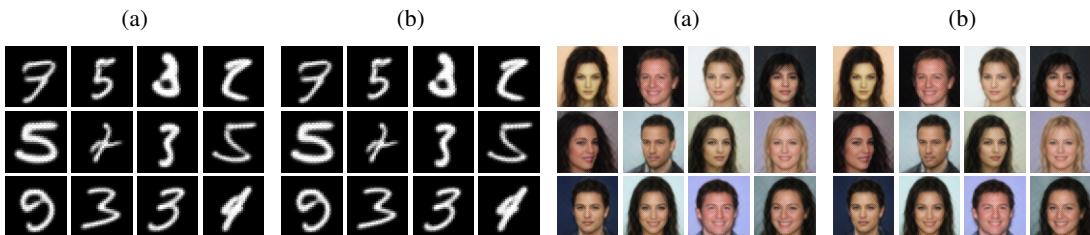


Figure 2: Comparison between multi-step and one-step flow matching. Panel labels: **(a)** multi-step Linear FM. **(b)** one-step Linear flow matching (FM);

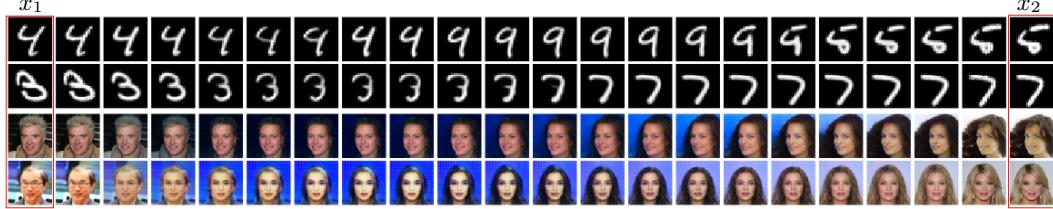


Figure 3: **One-step inversion examples:** Left and right (in red): original (not generated) data x_1 and x_2 . Intermediate images obtained by latent interpolation. See Appendix E.2 for more and higher-resolution results.

Geometric intuition. g acts as a diffeomorphism: in its latent space, the path between x_0 and x_1 is a straight line, and the velocity is constant. When mapped back to the data space via g^{-1} , this straight line becomes a curved trajectory. The Linearizer exploits this change of coordinates so that, in the right space, the velocity field is trivial, and the expensive integration disappears.

Implementation. The operator A is produced by an MLP whose input is t . To avoid huge matrices multiplications we build A as a low-rank matrix, by multiplying two rectangular matrices. g does not take t as input. Full implementation details are provided in App. G.

One-step generation. Figure 2 shows qualitative results on MNIST LeCun et al. (1998) (32×32) and CelebA Liu et al. (2015) (64×64). As discussed, our framework supports both multi-step and one-step sampling and yields effectively identical outputs. Visually, the samples are indistinguishable across datasets; quantitatively, the mean squared error between one-step and multi-step generations is 3.0×10^{-4} , confirming their near-equivalence in both datasets. Additionally, we evaluate FID across different step counts and our one-step simulation (see Figure 4(a)). The FID from 100 full iterative steps closely matches our one-step formulation, empirically validating the method’s theoretical guarantee of one-step sampling. Moreover, simulating more steps ($1000 \rightarrow 1$ vs. $100 \rightarrow 1$) improves FID by ~ 8 points, highlighting the strength of the linear-operator formulation. Furthermore, we validate one-step vs. 100 step fidelity, showing high similarity as presented in Figure 4(c). Finally, while our absolute FID is not yet competitive with state-of-the-art systems, our goal here is to demonstrate the theory in practice rather than to exhaustively engineer for peak performance; scaling and optimization are left to future work.

Inversion and interpolation. A fundamental limitation of flow models is that, in contrast to VAEs (Kingma & Welling, 2014), they lack a natural encoder: they cannot map data back into the prior (noise) space, and thus act only as decoders. As a result, inversion methods for diffusion (Dhariwal & Nichol, 2021; Song et al., 2021a; Huberman-Spiegelglas et al., 2024) have become an active research area. However, these techniques are approximate and often suffer from reconstruction errors, nonuniqueness, or computational overhead. Our framework offers a new bridge between diffusion-based models and encoding ability.

Because the Linearizer is linear, its Moore–Penrose pseudoinverse is also a Linearizer:

Lemma 7 (Moore–Penrose pseudoinverse of a Linearizer). *The Moore–Penrose pseudoinverse of f with respect to the induced inner products is*

$$f^\dagger(y) = g_x^{-1}(A^\dagger g_y(y)) \quad (19)$$

Proof. We verify the four Penrose equations (Penrose, 1955) in Appendix D. \square

This property enables the exact encoding of data in the latent space. For example, given two data points x_1, x_2 , we can encode them via $z^a = (1 - a) f^\dagger(x_1) + a f^\dagger(x_2)$, and decode back by

Full Steps			One-steps		PSNR		LPIPS	
1 10 100			(100→1) (1000→1)		MNIST	31.6	.008	MNIST
405	152	127	131	123	CelebA	33.4	.006	CelebA

(a) FID comparison.

(b) Inverse reconstruct consistency. (c) 100 vs. 1 step fidelity.

Figure 4: Quantitative comparisons.



Figure 5: **Style transfer examples.** Left: original image. Middle: style transfer using the left-side and right-side style images. Right: interpolation between the two styles.

$\hat{x}^a = f(z^a)$. Fig 3 shows latent interpolation between two real (non-generated) images. Additionally, we evaluate reconstruction quality using two standard metrics: LPIPS Zhang et al. (2018) and PSNR. Figure 4(b) shows inversion-reconstruction consistency. our results confirm that the proposed approach preserves information with high quality.

3.2 MODULAR STYLE TRANSFER

Style transfer has been widely studied since Gatys et al. (2016) proposed optimizing image pixels to match style and content statistics. Johnson et al. (2016) introduced perceptual-loss training of feed-forward networks, making style transfer practical.

Linearizer formulation. We fix g_x and g_y and associate each style with a matrix A_{style} . Analogously to A_t , the matrix A_{style} is produced by a hypernetwork that takes a *style index* as input (rather than time, as in the one-step setting). Then

$$f_{\text{style}}(x) = g_y^{-1}(A_{\text{style}}g_x(x)). \quad (20)$$

This separates content representation from style, making styles modular and algebraically manipulable. In practice, we distill A_{style} from a pretrained Johnson-style network.

Style interpolation. Given two trained style operators, A_{style}^x and A_{style}^y , we form an interpolated operator $A_{\text{style}}^{(\alpha x + (1-\alpha)y)}$ that represents a linear interpolation between them. We evaluated $\alpha \in \{0, 0.35, 0.40, 0.45, 0.50, 0.55, 0.60, 0.65, 1\}$ and computed

$$f_{\text{style}}^{(\alpha x + (1-\alpha)y)}(x) = g_y^{-1}\left(A_{\text{style}}^{(\alpha x + (1-\alpha)y)}g_x(x)\right),$$

with results shown in Figure 5. In the first row, we interpolate between the *mosaic* and *candy* styles; in the second row between *rain-princess* and *udine*; and in the third row between *candy* and *rain-princess*. See App.E for additional transfers.

3.3 LINEAR IDEMPOTENT GENERATIVE NETWORK

Idempotency, $f(f(x)) = f(x)$, is a central concept in algebra and functional analysis. In machine learning, it has been used in Idempotent Generative Networks (IGNs) (Shocher et al., 2024), to create a projective generative model. It was also used for robust test time training (Durasov et al.,

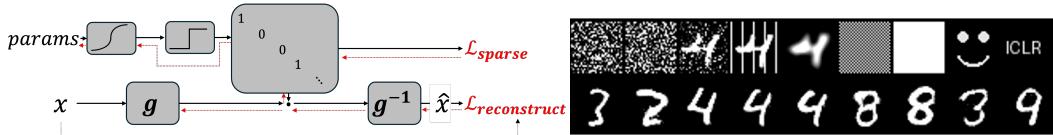


Figure 6: **Left: IGN training diagram.** Black solid arrows denote the forward pass; the red dashed arrow shows backpropagation. The parameter logits are passed through a sigmoid and then thresholded to form the binary projector A ; during backprop, the straight-through estimator (STE) bypasses the threshold so gradients flow through the underlying probabilities. **Right: Projection results.** Our Linear IGN makes a global projector that projects any input to the target distribution. Top are inputs and bottom are matching outputs.

2025). Enforcing a network to be idempotent is tricky and is currently done by using sophisticated optimization methods such as (Jensen & Vicary, 2025). The result is only an approximately idempotent model over the training data. We demonstrate enforcing accurate idempotency through architecture using our Linearizer. The key observation is that idempotency is preserved between the inner matrix A and the full function f .

Lemma 8. *The function f is idempotent \iff The matrix A is idempotent.*

Proof. For all x ,

$$\begin{aligned} A^2 = A &\iff g^{-1}(A^2 g(x)) = g^{-1}(Ag(x)) \\ &\iff \{g^{-1} \circ A \circ \underbrace{(g \circ g^{-1})}_{\text{Id}} \circ A \circ g\}(x) = \{g^{-1} \circ A \circ g\}(x) \\ &\iff f(f(x)) = f(x) \quad (\text{by the definition of } f) \quad \square \end{aligned}$$

Method. Figure 6 left shows the method. Recall that idempotent (projection) matrices have eigenvalues that are either 0 or 1. We assume a diagonalizable projection matrix $A = Q\Lambda Q^{-1}$ where Λ is a diagonal matrix with entries $\{0, 1\}$. The matrices Q, Q^{-1} can be absorbed into g, g^{-1} without loss of expressivity. So we can train a Linearizer $g^{-1}(\Lambda g(\cdot))$. In order to have Λ that is binary yet differentiable, we use straight-through estimation (Bengio et al., 2013) having underlying probabilities P as parameters in $[0, 1]$ and then use $A = \text{round}(P) + P - P.\text{detach}()$ where detach means stopping gradients. We forward propagate rounded values (0 or 1) and back propagate continuous values.

Idempotency by architecture allows reducing the idempotent loss used in Shocher et al. (2024). We need the data to be fixed points with the tightest possible latent manifold. The losses are:

$$\mathcal{L}_{rec} = \|f(x) - x\|^2 \tag{21}$$

$$\mathcal{L}_{sparse} = \frac{1}{N} \text{Rank}(A) = \frac{1}{N} \sum_i^N \lambda_i \tag{22}$$

We further enforce regularization that encourages g to preserve the norms.

$$\mathcal{L}_{isometry} = \left\| \|g(x) - g(0)\|^2 - \|x\|^2 \right\|_1 \tag{23}$$

This nudges g towards a near-isometry around the data, thus resisting mapping far-apart inputs to close outputs. It mitigates collapse and improves stability. We apply it with a small weight (0.001).

From local to global projectors. As Shocher et al. (2024) put it, “*We view this work as a first step towards a global projector.*” In practice, IGNs achieve idempotency only around the training data: near the source distribution or near the target distribution they are trained to reproduce. If you venture further away, then the outputs degenerate into arbitrary artifacts not related to the data distribution. In contrast, our Linearizer is idempotent by architecture. It does not need to be trained to approximate projection, it is one. Figure 6 right shows various inputs projected by our model onto the distribution. Interestingly, there is not even a notion of a separate *source distribution*: we never inject latent noise during training, and in a precise sense the entire ambient space serves as the source. This makes our construction a particularly unusual generative model and a natural next step toward the global projector envisioned in Shocher et al. (2024).

4 ANALYSIS AND DISCUSSION

4.1 EXPRESSIVENESS AND LIMITATIONS

A natural question is whether any arbitrary function can be represented by a Linearizer. The answer is no; the architecture imposes certain topological constraints. For example, the null space of f is completely determined by the null space of the matrix A . As linear operators can only have null spaces containing a single point (the zero vector) or an infinite subspace, f is restricted to these behaviors. This prevents a Linearizer from, for example, mapping three distinct points to zero.

However, for many practical machine learning tasks, we empirically find that this theoretical limitation is less restrictive than it appears. High-dimensional data, such as images or text embeddings,

often reside on a low-dimensional manifold granting the network a large degree of freedom for mapping off-manifold inputs. We hypothesize that in such cases most of the expressiveness constraints may not impact the task for the relevant data. Specifically, regarding the null space example, a finite number of on-manifold inputs can be mapped to zero, while infinitely many more off-manifold.

4.2 ON THE ROLE OF THE CORE OPERATOR A

One might also ask whether the matrix A is strictly necessary or if its action could be fully absorbed by the powerful invertible networks g_x and g_y . Let the SVD of A be $U\Sigma V^T$. In the general case where g_x and g_y are distinct, we could define new invertible networks $g'_y = g_y \circ U$ and $g'_x = V^T \circ g_x$. The function would then be $f(x) = (g'_y)^{-1}(\Sigma g'_x(x))$. If we further allow the networks to absorb the scaling defined by the singular values in Σ , then the core operator could be reduced to a diagonal matrix of zeros and ones (note that this does not imply idempotency unless $g_x = g_y$). In this sense, the fundamental role of A is simply to define the rank of the function f .

The situation is different in the constrained case where $g_x = g_y = g$. Here, the two networks are not independent. To preserve the structure $f(x) = g^{-1}(Ag(x))$, any transformation absorbed by g must be inverted by g^{-1} . This rules out g absorbing SVD elements U, V as they are not necessarily mutual inverses. For diagonalizable matrices however, Eigen-decomposition can be applied. Then, $A = Q\Lambda Q^{-1}$, where one could define $g' = Q^{-1} \circ g$ and $f(x) = (g')^{-1}(\Lambda g'(x))$. In this scenario, A cannot be reduced to a projection matrix, its role is to determine the eigenvalues of the function, a much stronger responsibility than just its rank. However, this is limited to diagonalizable matrices.

In practice, parameterizing an explicit and expressive A is needed in some cases, particularly when we wish to represent a family of functions (e.g., $A_1, A_2 \dots$) that share the same linear basis (g_x, g_y). The expressiveness of A allows more degrees of freedom for the diversity of different functions. For example, in our diffusion application, we learn a time-dependent operator A_t . A single pair of networks, g_x and g_y , defines a fixed basis, while explicit A_t allows for a rich variety of linear operations within that basis, controlled by the continuous parameter t .

5 LIMITATIONS

While the Linearizer introduces a principled framework for exact linearization of neural maps, several limitations remain. First, invertible networks are inherently more challenging to train than standard architectures, often requiring careful design. Second, this work represents a broad and general first step: our applications demonstrate feasibility, but none are yet scaled to state-of-the-art benchmarks. Finally, the precise expressivity of the Linearizer remains an open theoretical question.

6 RELATED WORK

Linearization in Dynamical Systems: The Koopman operator theory linearizes nonlinear dynamics by advancing observables on a single state space, yielding linear evolution $z_{t+1} = K z_t$ in discrete time or $\dot{z} = Az$ in continuous time (Mezić, 2005). Data-driven realizations such as DMD and EDMD construct finite-dimensional approximations of this infinite-dimensional operator on a chosen dictionary of observables, which typically produces a square matrix because the dictionary is mapped into itself, although variants with different feature sets exist (Schmid, 2010; Lusch et al., 2018; Mardt et al., 2018). Linear algebraic analyses such as SVD, eigendecomposition, and taking powers are standard in these Koopman approximations. Our work addresses a different object: an arbitrary learned mapping $f : \mathcal{X} \rightarrow \mathcal{Y}$. We learn invertible coordinate maps g_x and g_y and a finite matrix A such that $g_y \circ f \circ g_x^{-1} = A$, which gives the exact finite-dimensional linearity of f between induced spaces. This formulation allows for distinct input and output coordinates, allows A to be rectangular when $\dim \mathcal{X} \neq \dim \mathcal{Y}$, and transfers linear-algebraic structure in A directly to f in a controlled way. In particular, one can compute $\text{SVD}(f)$ via $\text{SVD}(A)$, perform analytic test-time inversion with A^+ , impose spectral constraints such as idempotency by restricting $\sigma(A)$, and more.

Neural Tangent Kernel: The NTK framework shows that infinitely wide networks trained with small steps evolve linearly in parameter space (Jacot et al., 2018). However, the resulting model remains nonlinear in input, output mapping. In comparison, our contribution achieves input-output linearity in a learned basis for any network width, independent of training dynamics.

One-Step Diffusion Distillation: Reducing sampling cost in diffusion models often uses student distillation: Progressive Distillation (Salimans & Ho, 2022), Consistency Models (Song et al., 2023), Distribution Matching Distillation (DMD) and f-distillation fall into this category. A recent work, the Koopman Distillation Model (KDM) (Berman et al., 2025), uses Koopman-based encoding to distill a diffusion model into a one-step generator, achieving strong FID improvements through offline distillation. The Koopman Distillation Model (KDM) distills a pre-trained diffusion teacher into an encoder linear decoder for diffusion only. We train from scratch, enforce exact finite-dimensional linearity by construction, and target multiple use cases including composition, inversion, idempotency, and continuous evolution. Unlike some invertible Koopman setups that force a shared map, KDM employs distinct encoder components and a decoder and does not enforce $g_x = g_y$.

Invertible Neural Networks: Invertible models like NICE (Dinh et al., 2015), RealNVP (Dinh et al., 2017) and Glow (Kingma & Dhariwal, 2018) use invertible transforms for density modeling and generative sampling. In this work, we make use of such invertible neural networks, to impose a bijective change of coordinates that allows linearity (lemma 2 only holds given this invertability of g_x, g_y) rather than modeling distribution.

Group-theoretic and Equivariant Representations: Let G act on \mathcal{X} via $(g, x) \mapsto g \cdot x$. Many works learn an encoder $e : \mathcal{X} \rightarrow V$ and a linear representation $\rho : G \rightarrow \text{GL}(V)$ such that the equivariance constraint $e(g \cdot x) = \rho(g)e(x)$ holds approximately, often with a decoder d that reconstructs x or enforces $d(\rho(g)z) \approx g \cdot d(z)$ (Quessard et al., 2020). The objective is to make the *action* of a symmetry group linear in latent coordinates. In contrast, we target a different object: an *arbitrary learned map* $f : \mathcal{X} \rightarrow \mathcal{Y}$. We construct invertible g_x, g_y and a finite matrix A so that $g_y \circ f \circ g_x^{-1} = A$, which yields *exact* finite-dimensional linearity of f between induced spaces and enables direct linear algebra on f via A .

Manifold Flattening and Diffeomorphic Autoencoders: Manifold-flattening methods assume data lie on a manifold $M \subset \mathbb{R}^N$ and learn or construct a map $\Phi : \mathbb{R}^N \rightarrow \mathbb{R}^k$ so that $\Phi(M)$ is close to a linear subspace L , often with approximate isometry on M (Psenka et al., 2024). Diffeomorphic autoencoders parameterize deformations $\varphi \in \text{Diff}(\Omega)$ with a latent z in a Lie algebra and use a decoder to warp a template, with variants using a log map to linearize the deformation composition (Bône et al., 2019). These approaches *flatten data geometry or deformation fields*. This contrasts with our Linearizer in which the flattening is over a mapping function rather than a dataset manifold or a deformation group.

Deep Linear Networks: Networks composed solely of linear layers are linear in the Euclidean basis and are primarily used to study optimization paths (Saxe et al., 2014; Arora et al., 2018; Cohen et al., 2016). They lack expressive power in standard coordinates. In contrast, our Linearizer is expressive thanks to g_x and g_y that are nonlinear over the standard vector spaces, yet maintains exact linearity in the induced coordinate system.

Normalizing flows: Normalizing flows Rezende & Mohamed (2015) use invertible neural networks to transport probability measures: they map a simple Gaussian distribution into a complex data distribution while preserving density through the change-of-variables formula. Linearizers use the same building block, an invertible network, but instead of transporting measures, they transport algebraic structure. In this sense, Linearizers may be viewed as the “linear-algebraic analogue” of normalizing flows: both frameworks exploit invertible transformations to pull back complex objects into domains where they become simple and tractable.

7 CONCLUSION

We introduced the Linearizer, a framework that learns invertible coordinate transformations such that neural networks become exact linear operators in the induced space. This yields new applications in one-step flow matching, modular style transfer, and architectural enforcement of idempotency.

Looking ahead, several directions stand out. First, scaling one-step flow matching and IGN to larger datasets and higher resolutions could provide competitive generative models with unprecedented efficiency. Second, the matrix structure of the Linearizer opens the door to modeling motion dynamics: by exploiting matrix exponentials, one could simulate continuous-time evolution directly, extending the approach beyond generation to physical and temporal modeling. More broadly, characterizing the theoretical limits of Linearizer expressivity, and integrating it with other operator-learning frameworks, remains an exciting avenue for future research.

ACKNOWLEDGEMENTS

We thank Amil Dravid and Yoad Tewel for insightful discussions and for sharing ideas that helped shape this work. A.S. is supported by the Chaya Career Advancement Chair.

REFERENCES

- Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit acceleration by overparameterization. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*, pp. 244–253. PMLR, 2018.
- Yoshua Bengio, Nicolas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons. In *Proceedings of the 2nd International Conference on Learning Representations (ICLR, Workshop Track)*, 2013.
- Nimrod Berman, Ilan Naiman, Moshe Eliasof, Hedi Zisling, and Omri Azencot. One-step offline distillation of diffusion-based models via koopman modeling. In *Proceedings of the 42nd International Conference on Machine Learning (ICML)*. PMLR, 2025.
- Alexandre Bône, Maxime Louis, Olivier Colliot, and Stanley Durrleman. Learning low-dimensional representations of shape data sets with diffeomorphic autoencoders. In *Information Processing in Medical Imaging (IPMI)*, volume 11492 of *Lecture Notes in Computer Science*, pp. 195–207. Springer, 2019.
- Nadav Cohen, Or Sharir, and Amnon Shashua. On the expressive power of deep learning: A tensor analysis. In *Proceedings of the 29th Annual Conference on Learning Theory (COLT)*, 2016.
- Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 34, 2021.
- Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Non-linear independent components estimation. In *Proceedings of the 2nd International Conference on Learning Representations (ICLR, Workshop Track)*, 2015.
- Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. In *Proceedings of the 5th International Conference on Learning Representations (ICLR)*, 2017.
- Nikita Durasov, Assaf Shocher, Doruk Oner, Gal Chechik, Alexei A. Efros, and Pascal Fua. IT³: Idempotent test-time training. In *Proceedings of the 42nd International Conference on Machine Learning (ICML)*. PMLR, 2025.
- Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2414–2423, 2016.
- Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5):359–366, 1989.
- Yarden Huberman-Spiegelglas, Shai Bagon, and Michal Irani. Diffusion inversion: A universal technique for diffusion editing and interpolation. *Transactions on Machine Learning Research (TMLR)*, 2024.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 31, 2018.
- Nikolaj Banke Jensen and Jamie Vicary. Enforcing idempotency in neural networks. In *Proceedings of the 42nd International Conference on Machine Learning (ICML)*, 2025.
- Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision (ECCV)*, pp. 694–711, 2016.

- Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *Advances in neural information processing systems*, 35:26565–26577, 2022.
- Diederik P. Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 31, 2018.
- Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *Proceedings of the 2nd International Conference on Learning Representations (ICLR)*, 2014.
- Wilhelm Kutta. *Beitrag zur näherungsweisen Integration totaler Differentialgleichungen*. Teubner, 1901.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. In *International Conference on Learning Representations (ICLR)*, 2023.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- Bethany Lusch, J. Nathan Kutz, and Steven L. Brunton. Deep learning for universal linear embeddings of nonlinear dynamics. *Nature Communications*, 9:4950, 2018.
- Andreas Mardt, Luca Pasquali, Hao Wu, and Frank Noé. Vampnets for deep learning of molecular kinetics. *Nature Communications*, 9:5, 2018.
- Igor Mezić. Spectral properties of dynamical systems, model reduction and decompositions. *Nonlinear Dynamics*, 41:309–325, 2005.
- Roger Penrose. A generalized inverse for matrices. *Mathematical Proceedings of the Cambridge Philosophical Society*, 51(3):406–413, 1955.
- Michael Psenka, Druv Pai, Vishal Raman, Shankar Sastry, and Yi Ma. Representation learning via manifold flattening and reconstruction. *Journal of Machine Learning Research*, 2024.
- Robin Quessard, Thomas D. Barrett, and William R. Clements. Learning group structure and disentangled representations of dynamical environments. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, 2020.
- Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, pp. 1530–1538. PMLR, 2015.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pp. 234–241. Springer, 2015.
- Carl Runge. Über die numerische auflösung von differentialgleichungen. *Mathematische Annalen*, 46(2):167–178, 1895.
- Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 35, 2022.
- Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. In *International Conference on Learning Representations (ICLR)*, 2014.
- Peter J. Schmid. Dynamic mode decomposition of numerical and experimental data. *Journal of Fluid Mechanics*, 656:5–28, 2010.
- Assaf Shocher, Amil David, Yossi Gandelsman, Inbar Mosseri, Michael Rubinstein, and Alexei A. Efros. Idempotent generative network. In *International Conference on Learning Representations (ICLR)*, 2024.

Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations (ICLR)*, 2021a.

Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations (ICLR)*, 2021b.

Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. In *Proceedings of the 40th International Conference on Machine Learning (ICML)*, 2023.

Gilbert Strang. *Introduction to Linear Algebra*. SIAM, 2022.

Steven H. Strogatz. *Nonlinear Dynamics and Chaos: With Applications to Physics, Biology, Chemistry, and Engineering*. Chapman and Hall/CRC, 2024.

Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 586–595, 2018.

A PROOF OF LEMMA 1 (VALID VECTOR SPACES)

Proof. We verify that (V, \oplus_g, \odot_g) is a vector space over \mathbb{R} by transporting each axiom via g and g^{-1} , using only the definitions

$$u \oplus_g v := g^{-1}(g(u) + g(v)), \quad a \odot_g u := g^{-1}(a g(u)).$$

For $u, v, w \in V$ and $a, b \in \mathbb{R}$:

1. Closure

$$u \oplus_g v = g^{-1}(g(u) + g(v)) \in V. \quad (24)$$

2. Associativity

$$\begin{aligned} (u \oplus_g v) \oplus_g w &= g^{-1}(g(u \oplus_g v) + g(w)) \\ &= g^{-1}((g(u) + g(v)) + g(w)) \\ &= g^{-1}(g(u) + (g(v) + g(w))) \\ &= g^{-1}(g(u) + g(v \oplus_g w)) \\ &= u \oplus_g (v \oplus_g w). \end{aligned} \quad (25)$$

3. Commutativity

$$\begin{aligned} u \oplus_g v &= g^{-1}(g(u) + g(v)) \\ &= g^{-1}(g(v) + g(u)) \\ &= v \oplus_g u. \end{aligned} \quad (26)$$

4. Additive identity (with $0_V := g^{-1}(0)$)

$$u \oplus_g 0_V = g^{-1}(g(u) + g(0_V)) = g^{-1}(g(u) + 0) = u. \quad (27)$$

5. Additive inverse (with $(-u) := g^{-1}(-g(u))$)

$$u \oplus_g (-u) = g^{-1}(g(u) + g(-u)) = g^{-1}(g(u) - g(u)) = g^{-1}(0) = 0_V. \quad (28)$$

6. Compatibility of scalar multiplication

$$\begin{aligned} a \odot_g (b \odot_g u) &= g^{-1}(a g(b \odot_g u)) = g^{-1}(a(b g(u))) \\ &= g^{-1}((ab) g(u)) = (ab) \odot_g u. \end{aligned} \quad (29)$$

7. Scalar identity

$$1 \odot_g u = g^{-1}(1 \cdot g(u)) = g^{-1}(g(u)) = u. \quad (30)$$

8. Distributivity over vector addition

$$\begin{aligned} a \odot_g (u \oplus_g v) &= g^{-1}(a g(u \oplus_g v)) = g^{-1}(a(g(u) + g(v))) \\ &= g^{-1}(a g(u) + a g(v)) = g^{-1}(a g(u)) \oplus_g g^{-1}(a g(v)) \\ &= (a \odot_g u) \oplus_g (a \odot_g v). \end{aligned} \quad (31)$$

9. Distributivity over scalar addition

$$\begin{aligned} (a + b) \odot_g u &= g^{-1}((a + b) g(u)) = g^{-1}(a g(u) + b g(u)) \\ &= g^{-1}(a g(u)) \oplus_g g^{-1}(b g(u)) = (a \odot_g u) \oplus_g (b \odot_g u). \end{aligned} \quad (32)$$

□

B PROOF OF LEMMA 6 (SVD OF A LINEARIZER)

Let $f : \mathcal{X} \rightarrow \mathcal{Y}$ be a Linearizer

$$f(x) = g_y^{-1}(A g_x(x)), \quad (33)$$

where $g_x : \mathcal{X} \rightarrow \mathbb{R}^N$ and $g_y : \mathcal{Y} \rightarrow \mathbb{R}^M$ are invertible and $A \in \mathbb{R}^{M \times N}$. Equip \mathcal{X} and \mathcal{Y} with the induced inner products

$$\langle u, v \rangle_{g_x} := \langle g_x(u), g_x(v) \rangle_{\mathbb{R}^N}, \quad \langle p, q \rangle_{g_y} := \langle g_y(p), g_y(q) \rangle_{\mathbb{R}^M}. \quad (34)$$

Let the (Euclidean) SVD of A be $A = U\Sigma V^\top$, where $U = [u_1, \dots, u_M] \in \mathbb{R}^{M \times M}$, $V = [v_1, \dots, v_N] \in \mathbb{R}^{N \times N}$ are orthogonal and $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots)$ with $\sigma_1 \geq \dots \geq \sigma_r > 0$. Define

$$\tilde{u}_i := g_y^{-1}(u_i) \in \mathcal{Y}, \quad \tilde{v}_i := g_x^{-1}(v_i) \in \mathcal{X}. \quad (35)$$

Then $\{\tilde{v}_i\}$ and $\{\tilde{u}_i\}$ are orthonormal sets in $(\mathcal{X}, \langle \cdot, \cdot \rangle_{g_x})$ and $(\mathcal{Y}, \langle \cdot, \cdot \rangle_{g_y})$, respectively, and for each $i \leq r$,

$$f(\tilde{v}_i) = \sigma_i \odot_{g_y} \tilde{u}_i, \quad \text{and} \quad f^*(\tilde{u}_i) = \sigma_i \odot_{g_x} \tilde{v}_i, \quad (36)$$

where f^* is the adjoint with respect to the induced inner products and $a \odot_g$ denotes the induced scalar multiplication. Hence $(\{\tilde{u}_i\}_{i=1}^r, \{\sigma_i\}_{i=1}^r, \{\tilde{v}_i\}_{i=1}^r)$ is an SVD of f between the induced Hilbert spaces.

Proof. Orthonormality. For i, j ,

$$\langle \tilde{v}_i, \tilde{v}_j \rangle_{g_x} = \langle g_x(\tilde{v}_i), g_x(\tilde{v}_j) \rangle_{\mathbb{R}^N} = \langle v_i, v_j \rangle_{\mathbb{R}^N} = \delta_{ij}, \quad (37)$$

and similarly $\langle \tilde{u}_i, \tilde{u}_j \rangle_{g_y} = \delta_{ij}$. Thus g_x and g_y are isometric isomorphisms from the induced spaces to Euclidean space.

Action on right singular vectors. Using $Av_i = \sigma_i u_i$,

$$f(\tilde{v}_i) = g_y^{-1}(A g_x(\tilde{v}_i)) = g_y^{-1}(Av_i) = g_y^{-1}(\sigma_i u_i) = \sigma_i \odot_{g_y} \tilde{u}_i. \quad (38)$$

Adjoint and action on left singular vectors. The adjoint $f^* : \mathcal{Y} \rightarrow \mathcal{X}$ with respect to the induced inner products is characterized by

$$\langle f(x), y \rangle_{g_y} = \langle x, f^*(y) \rangle_{g_x} \quad \text{for all } x \in \mathcal{X}, y \in \mathcal{Y}. \quad (39)$$

Transporting through g_x, g_y and using Euclidean adjoints shows that

$$f^*(y) = g_x^{-1}(A^\top g_y(y)). \quad (40)$$

Therefore, since $A^\top u_i = \sigma_i v_i$,

$$f^*(\tilde{u}_i) = g_x^{-1}(A^\top g_y(\tilde{u}_i)) = g_x^{-1}(A^\top u_i) = g_x^{-1}(\sigma_i v_i) = \sigma_i \odot_{g_x} \tilde{v}_i. \quad (41)$$

Conclusion. The triples $(\tilde{u}_i, \sigma_i, \tilde{v}_i)$ satisfy the defining relations of singular triplets for f between the induced inner product spaces, with $\{\tilde{v}_i\}$ and $\{\tilde{u}_i\}$ orthonormal bases of the right and left singular subspaces. For zero singular values, the same construction holds with images mapped to the null spaces as usual. \square

C PROOF OF LEMMA 4 (HILBERT SPACE)

Let $g : V \rightarrow \mathbb{R}^n$ be a bijection and endow V with the induced vector space operations (\oplus_g, \odot_g) and inner product

$$\langle u, v \rangle_g := \langle g(u), g(v) \rangle_{\mathbb{R}^n}. \quad (42)$$

Then $(V, \langle \cdot, \cdot \rangle_g)$ is a Hilbert space.

Proof. By Lemma 1, (V, \oplus_g, \odot_g) is a vector space and g is a vector-space isomorphism from (V, \oplus_g, \odot_g) onto $(\mathbb{R}^n, +, \cdot)$. The form $\langle \cdot, \cdot \rangle_g$ is an inner product because it is the pullback of the Euclidean inner product: bilinearity, symmetry, and positive definiteness follow immediately from injectivity of g and the corresponding properties in \mathbb{R}^n .

Let $\|\cdot\|_g$ be the norm induced by $\langle \cdot, \cdot \rangle_g$. For any $u, v \in V$,

$$\|u - v\|_g = \|g(u) - g(v)\|_2, \quad (43)$$

so g is an isometry from $(V, \|\cdot\|_g)$ onto $(\mathbb{R}^n, \|\cdot\|_2)$. Hence a sequence $\{u_k\} \subset V$ is Cauchy in $\|\cdot\|_g$ iff $\{g(u_k)\}$ is Cauchy in \mathbb{R}^n . Since \mathbb{R}^n is complete, $\{g(u_k)\}$ converges to some $y \in \mathbb{R}^n$. By surjectivity of g , there exists $u^* \in V$ with $g(u^*) = y$, and then $\|u_k - u^*\|_g = \|g(u_k) - y\|_2 \rightarrow 0$. Thus $(V, \|\cdot\|_g)$ is complete. Therefore $(V, \langle \cdot, \cdot \rangle_g)$ is a Hilbert space. \square

D PROOF OF LEMMA 7 (PSEUDO-INVVERSE OF A LINEARIZER)

Pseudoinverse of a Linearizer. Let $f : \mathcal{X} \rightarrow \mathcal{Y}$ be a Linearizer $f(x) = g_y^{-1}(A g_x(x))$ under the induced inner products $\langle u, v \rangle_{g_x} := \langle g_x(u), g_x(v) \rangle$ on \mathcal{X} and $\langle u, v \rangle_{g_y} := \langle g_y(u), g_y(v) \rangle$ on \mathcal{Y} . Write $f^*(y) = g_x^{-1}(A^\top g_y(y))$ for the adjoint and let A^\dagger be the (Euclidean) Moore–Penrose pseudoinverse of A (Penrose, 1955).

The Moore–Penrose pseudoinverse of f with respect to the induced inner products is

$$f^\dagger = g_x^{-1} \circ A^\dagger \circ g_y. \quad (44)$$

Proof. Following Penrose (1955), the Moore–Penrose pseudoinverse of a linear operator is uniquely characterized as the map satisfying four algebraic conditions (the Penrose equations). To establish that $f^\dagger = g_x^{-1} \circ A^\dagger \circ g_y$ is indeed the pseudoinverse of f , it therefore suffices to verify these four identities explicitly:

1. $ff^\dagger f = f$,
2. $f^\dagger f f^\dagger = f^\dagger$,
3. $(ff^\dagger)^* = ff^\dagger$,
4. $(f^\dagger f)^* = f^\dagger f$.

We verify the Penrose equations in the induced spaces.

(1) $ff^\dagger f = f$:

$$\begin{aligned} f f^\dagger f &= (g_y^{-1} \circ A \circ g_x) (g_x^{-1} \circ A^\dagger \circ g_y) (g_y^{-1} \circ A \circ g_x) \\ &= g_y^{-1}(AA^\dagger A g_x) \\ &= g_y^{-1}(A g_x) = f, \end{aligned} \quad (45)$$

using $AA^\dagger A = A$.

(2) $f^\dagger f f^\dagger = f^\dagger$:

$$\begin{aligned} f^\dagger f f^\dagger &= (g_x^{-1} \circ A^\dagger \circ g_y) (g_y^{-1} \circ A \circ g_x) (g_x^{-1} \circ A^\dagger \circ g_y) \\ &= g_x^{-1}(A^\dagger AA^\dagger g_y) \\ &= g_x^{-1}(A^\dagger g_y) = f^\dagger, \end{aligned} \quad (46)$$

using $A^\dagger AA^\dagger = A^\dagger$.

(3) $(ff^\dagger)^* = ff^\dagger$ on \mathcal{Y} :

$$\begin{aligned} ff^\dagger &= g_y^{-1}(AA^\dagger g_y), \\ (ff^\dagger)^* &= g_y^{-1}((AA^\dagger)^\top g_y) = g_y^{-1}(AA^\dagger g_y) = ff^\dagger, \end{aligned} \quad (47)$$

since AA^\dagger is symmetric.

(4) $(f^\dagger f)^* = f^\dagger f$ on \mathcal{X} :

$$\begin{aligned} f^\dagger f &= g_x^{-1}(A^\dagger A g_x), \\ (f^\dagger f)^* &= g_x^{-1}((A^\dagger A)^\top g_x) = g_x^{-1}(A^\dagger A g_x) = f^\dagger f, \end{aligned} \quad (48)$$

since $A^\dagger A$ is symmetric. All four conditions hold, hence f^\dagger is the Moore–Penrose pseudoinverse of f in the induced inner-product spaces. \square

E ADDITIONAL RESULTS

E.1 STYLE TRANSFER

We expand the results of style transfer made in the main paper and show multiple more transformations in Figure 7 following the same setup.



Figure 7: Style transfer examples. Left: original image. Middle: style transfer using the left-side and right-side style images. Right: interpolation between the two styles.

E.2 ONE-STEP INVERSE INTERPOLATION

We extend the results presented in the main text by providing additional interpolation examples on the MNIST (Figure 9) and CelebA (Figure 8) datasets.



Figure 8: Additional interpolations of CelebA via inversion.

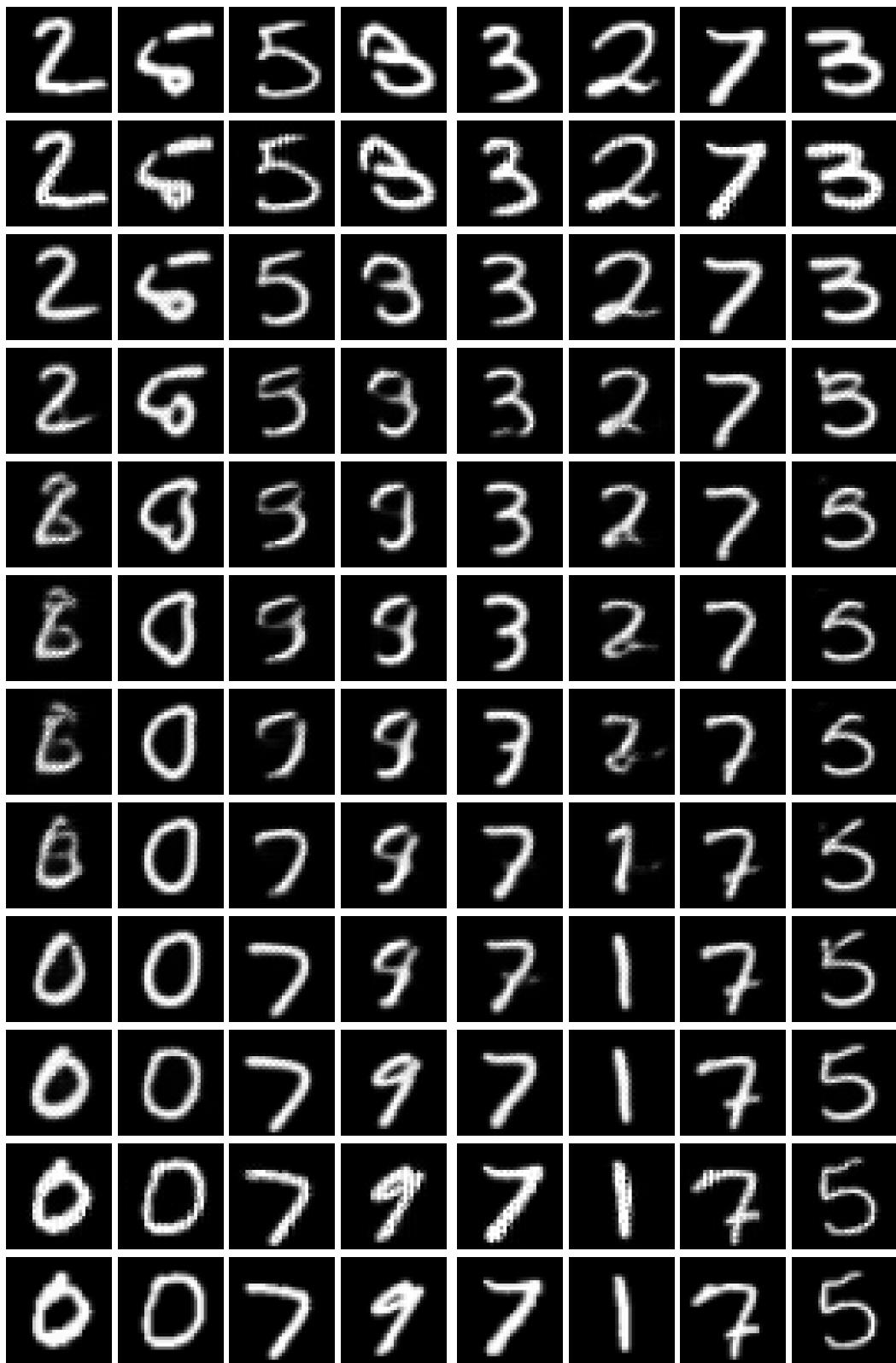


Figure 9: Additional interpolations of MNIST via inversion.

F DERIVATION FOR RUNGE-KUTTA COLLAPSE TO ONE STEP

We follow the setup in Sec. 3.1. Let $z_t := g(x_t)$ and $f(x, t) = g^{-1}(A_t g(x))$. In the induced coordinates, addition and scaling are Euclidean, so the ODE update acts linearly on z_t .

RK4 step in the induced space. Define $t_n = n\Delta t$ and $z_n := z_{t_n} = g(x_{t_n})$. The classical RK4 step from t_n to $t_{n+1} = t_n + \Delta t$ is

$$k_1 = A_{t_n} z_n, \quad (49)$$

$$k_2 = A_{t_n + \frac{\Delta t}{2}} \left(z_n + \frac{\Delta t}{2} k_1 \right), \quad (50)$$

$$k_3 = A_{t_n + \frac{\Delta t}{2}} \left(z_n + \frac{\Delta t}{2} k_2 \right), \quad (51)$$

$$k_4 = A_{t_n + \Delta t} \left(z_n + \Delta t k_3 \right), \quad (52)$$

$$z_{n+1} = z_n + \frac{\Delta t}{6} (k_1 + 2k_2 + 2k_3 + k_4). \quad (53)$$

Since each k_i is linear in z_n , there exists a matrix M_n such that $z_{n+1} = M_n z_n$. Expanding the linearity gives the explicit one-step operator

$$M_n = I + \frac{\Delta t}{6} \left[A_{t_n} + 2 A_{t_n + \frac{\Delta t}{2}} \left(I + \frac{\Delta t}{2} A_{t_n} \right) \right] \quad (54)$$

$$+ 2 A_{t_n + \frac{\Delta t}{2}} \left(I + \frac{\Delta t}{2} A_{t_n + \frac{\Delta t}{2}} \left(I + \frac{\Delta t}{2} A_{t_n} \right) \right) \quad (55)$$

$$+ A_{t_n + \Delta t} \left(I + \Delta t A_{t_n + \frac{\Delta t}{2}} \left(I + \frac{\Delta t}{2} A_{t_n + \frac{\Delta t}{2}} \left(I + \frac{\Delta t}{2} A_{t_n} \right) \right) \right)]. \quad (56)$$

Collapsed one-step operator. Iterating N steps and collecting the product in latent space yields

$$B := \prod_{n=0}^{N-1} M_n, \quad \text{so that} \quad g(\hat{x}_1) = B g(x_0), \quad \hat{x}_1 = g^{-1}(B g(x_0)). \quad (57)$$

The last two equations (equation 54 and equation 57) are the RK4 analogues of the Euler collapse in the main text: the multi-step solver is replaced by a single multiplication with B computed once after training.

Time-invariant sanity check (optional). If $A_t \equiv A$ is constant, then equation 54 reduces to the degree-4 Taylor polynomial of $e^{\Delta t A}$:

$$M_n = I + \Delta t A + \frac{(\Delta t)^2}{2} A^2 + \frac{(\Delta t)^3}{6} A^3 + \frac{(\Delta t)^4}{24} A^4, \quad B = M^N. \quad (58)$$

This matches the behavior of classical RK4 on linear time-invariant systems.

G IMPLEMENTATION DETAILS

G.1 BACKBONES AND BLOCK STRUCTURE

We use **6 invertible blocks** in all models. The two architectures differ only in the internal block type:

Model family	Per-block flow (forward)
Diffusion / Style	(optional Squeeze2x2 for 1ch) \rightarrow ActNorm \rightarrow Affine Coupling ($x_1 x_2$) \rightarrow Affine Coupling ($x_2 y_1$) \rightarrow Invertible 1×1 Conv
IGN	SpatialSplit2x2Rand \rightarrow Additive Couplings with $\{F, G\}$ \rightarrow concat \rightarrow PixelShuffle ($2 \times$)

Table 1: Block-level flow for the two architectures; repeated 6 times.

Notes.

- **ActNorm**: per-channel affine with data-dependent init on first batch.
- **Affine coupling (Diffusion/Style)**: shift/log-scale predicted by a tiny U-Net conditioner; clamped log s .
- **Invertible 1×1 conv (Diffusion/Style)**: learned channel permutation/mixer (orthogonal init).
- **SpatialSplit2x2Rand (IGN)**: deterministic 2×2 space partition into two streams; inverse merges.
- **PixelShuffle / unshuffle (IGN)**: reversible $2 \times$ spatial reindexing between blocks.
- **We do not use InvTanhScaled**.

The blocks:

1. **Diffusion/Style.** We use a U-Net (Ronneberger et al., 2015), adapting the implementation of Karras et al. (2022). We modify only two hyperparameters: `model_channels = 16` and `channel_mult = [1, 1]`.
2. **IGN.** We use bottleneck architecture, gradually downscaling to spatial size 1×1 by stride 2 convolutions, and then using transposed convolutions back to original resolution. See table below for exact layers.

G.2 TIME PATH (DIFFUSION)

In diffusion/flow-matching models, the scalar t is embedded and fed to: (i) the affine-coupling conditioners inside g , and (ii) the time-dependent core A_t (low-rank MLP). Style and IGN models do not use t .

G.3 CORES A

- **Diffusion:** Two non-linear MLPs that produce two matrices for the low-rank parameterization $A = A_1 A_2$ with $\text{rank}(A) = 16$.
- **Style.** A kernel 4×4 is produced by a non-linear hypernetwork $L(\text{style}) = A_{\text{ker}}^{\text{style}}$, which we then apply to the image via a 2D convolution (PyTorch).
- **IGN:** diagonal A (elementwise scale in latent); for idempotent IGN we threshold the diagonal with STE.

G.4 CONDITIONERS

Diffusion / Style (Affine-coupling conditioners). Each coupling uses a tiny U-Net that maps $C_{\text{in}} = C/2 \rightarrow 2 C_{\text{in}}$ (shift, log s).

#	Layer (per conditioner)
1	GroupNorm → Conv3 × 3 → SiLU
2	Residual block × 2 (GroupNorm, Conv3 × 3, SiLU)
3	Final Conv3 × 3 to 2 C_{in} ; clamp log s

Table 2: Minimal U-Net-style conditioner used in affine couplings.

IGN (Additive-coupling CNNs). Each coupling uses a plain CNN (*no U-Net skip concatenations*); forward is additive, inverse is exact.

#	F or G layer stack
1	Conv2d 2 → 8, kernel 4, stride 2
2	Conv2d 8 → 32, kernel 4, stride 2
3	Conv2d 32 → 128, kernel 4, stride 2
4	Conv2d 128 → 512, kernel 4, stride 2
5	ConvTranspose2d 512 → 128, kernel 4, stride 2
6	ConvTranspose2d 128 → 32, kernel 4, stride 2
7	ConvTranspose2d 32 → 8, kernel 4, stride 2
8	ConvTranspose2d 8 → 2, kernel 4, stride 2 → tanh

Table 3: CNN used for F and G inside IGN additive couplings.

G.5 LOSSES (WEIGHTS FIXED)

IGN.

$$\lambda_{\text{rec}} = 1.0, \quad \lambda_{\text{sparse}} = 0.75, \quad \lambda_{\text{iso}} = 0.001.$$

We use (i) reconstruction MSE, (ii) sparsity/rank on the diagonal of A , and (iii) a light isometry term on g .

Diffusion / Flow-Matching. We use the main training loss described in the paper. For additional stability, we empirically found that adding the perceptual reconstruction terms $d(x_0, g^{-1}(g(x_0)))$ and $d(x_1, g^{-1}(g(x_1)))$ (with d as LPIPS (Zhang et al., 2018)), together with the alignment term $\|A g(x_0) - g(x_1)\|^2$, yields more stable training and higher-quality models. The primary loss in the paper also uses LPIPS as the distance metric. No explicit weighting between losses was required.

G.6 ONE-LINE HYPERPARAMETER SUMMARY

Item	Setting
Invertible blocks	6 (both families)
Permutations	Diff/Style: invertible 1 × 1 conv; IGN: SpatialSplit2x2Rand + PixelShuffle
ActNorm	Enabled (both)
Core A	Diff/Style: low-rank 16; IGN: diagonal (STE for projector)
Time t	Used in Diffusion/Style (conditioners + A_t); not used in IGN

Table 4: Everything needed to reproduce the scaffolding.