

# ICASSP 2024 Tutorial T-10: Building White-Box Deep Neural Networks

---

*Lecture III: White-Box Autoencoding via Structured Denoising-Diffusion* **Sam Buchanan**, TTIC

**Druv Pai**, UC Berkeley

April 14, 2024



## Lecture III Summary

- Today so far:
  1. Unrolled optimization and foundations of white-box models
  2. The model problem of compressive autoencoding, and an objective which promotes compressive representations
  3. A white-box transformer encoder which performs compressive encoding at scale
- Now:
  4. A white-box transformer *autoencoder* for compressive autoencoding, built using connections between classical signal processing primitives and tasks
  5. Beyond autoencoding for representation learning; principled self-consistent closed-loop learning systems

# Outline

---

White-Box Autoencoders via Structured Denoising-Diffusion

Closed-Loop Transcription: Self-Consistent Representation Learning

Tutorial Conclusion

## Recall: Compressive Autoencoding

Abstractly: seek maps  $f : \mathbb{R}^D \rightarrow \mathbb{R}^d$ ,  $g : \mathbb{R}^d \rightarrow \mathbb{R}^D$ , such that

$$\mathbf{x} \in \mathbb{R}^D \xrightarrow{f(\mathbf{x})} \mathbf{z} \in \mathbb{R}^d \xrightarrow{g(\mathbf{z})} \hat{\mathbf{x}} \in \mathbb{R}^D.$$

satisfies  $\mathbf{x} \approx \hat{\mathbf{x}} = g \circ f(\mathbf{x})$  in a suitable sense; and  $\mathbf{z}$  is

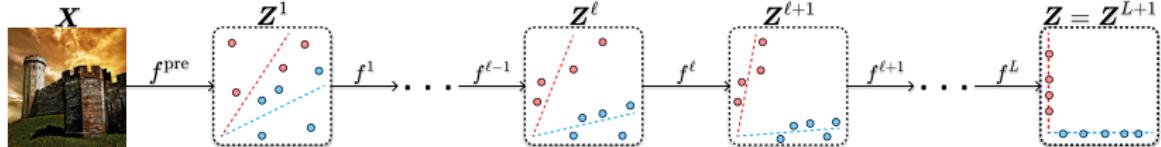
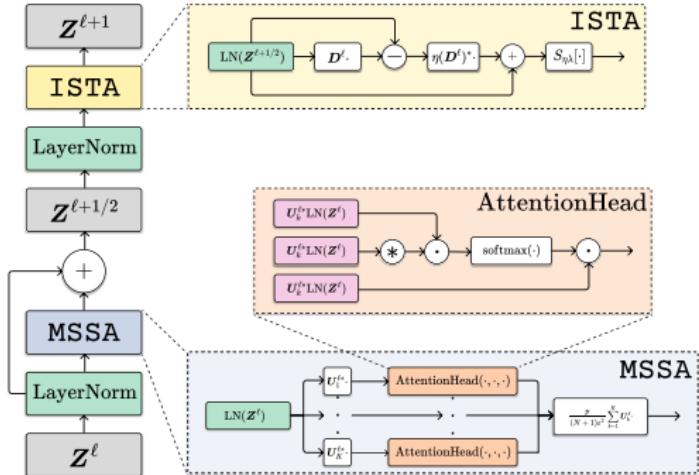
1. **Compressed**, preserving only intrinsic low-dimensionality of  $\mathbf{x}$ ;
2. **Linearized**, having (piecewise) linear geometry, to aid interpolation and extrapolation in the representation space;
3. **Sparse**: distinct modes of variability in  $\mathbf{x}$  are geometrically *incoherent* and moreover *axis-aligned*.

**Goal:** Build a decoder where we have the same level of layer-wise understanding as the encoder.

# White-Box Autoencoder

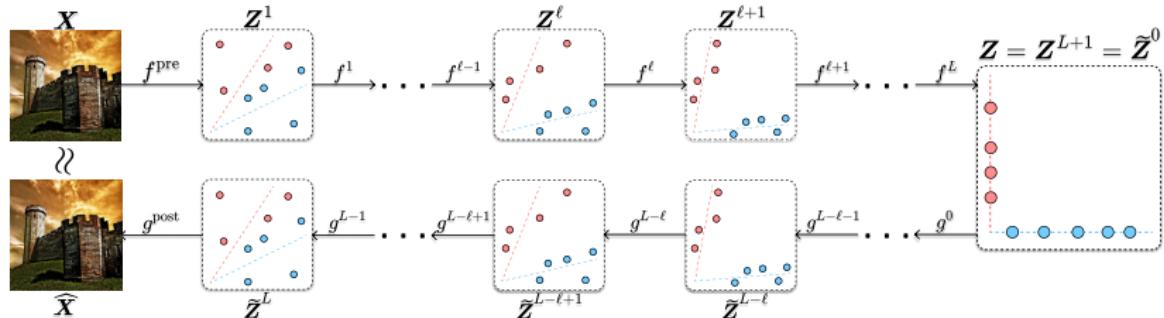
We understand *encoder* via unrolled optimization and compression.

What about the *decoder*?



# Unrolled Optimization for Building White-Box Decoder?

Decoder **cannot** optimize representation learning objective.

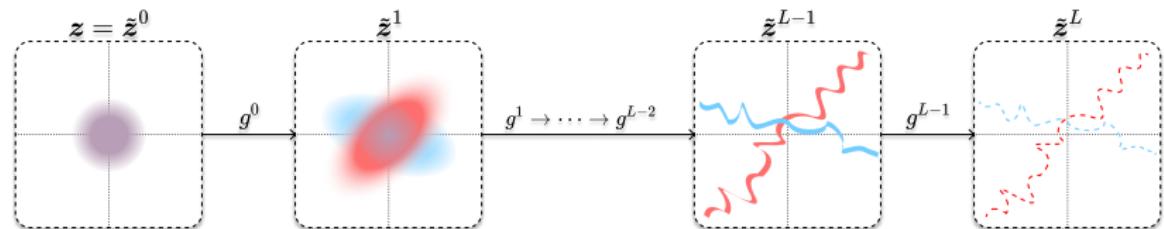


Want *reversible process* which converts *noisy, nonlinear* data to *compact, structured* representations (and back).

# Diffusion for Building White-Box Decoder?

**Diffusion process:** scalable, (distribution-)invertible way to convert:

random noise  $\longleftrightarrow$  low-dim. data distribution



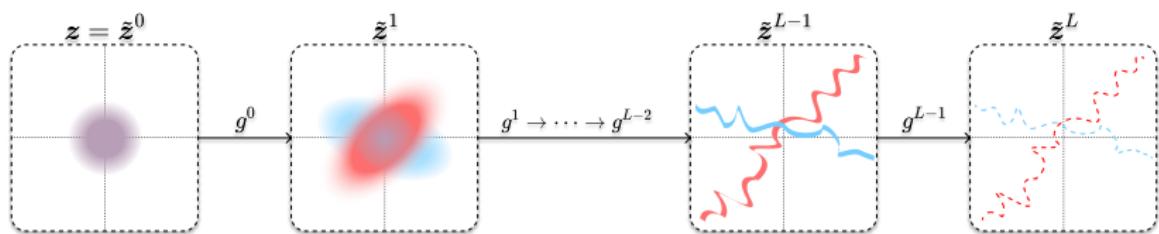
$$\tilde{z}^{\ell+1} = g^\ell(\tilde{z}^\ell) = \tilde{z}^\ell + \tau^\ell \nabla \log q^\ell(\tilde{z}^\ell)$$

# Denoising and Diffusion Models

**Tweedie's formula** [Robbins 1956]: If  $\mathbf{z}^\ell = \mathbf{z}_\natural^\ell + (\sigma^\ell)^2 \mathbf{g}^\ell$ , then

$$\mathbb{E}[\mathbf{z}_\natural^\ell | \mathbf{z}^\ell] = \mathbf{z}^\ell + (\sigma^\ell)^2 \nabla \log q^\ell(\mathbf{z}^\ell),$$

where  $q^\ell \doteq \text{PDF}(\mathbf{z}^\ell)$  and  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . C.f. diffusion models:



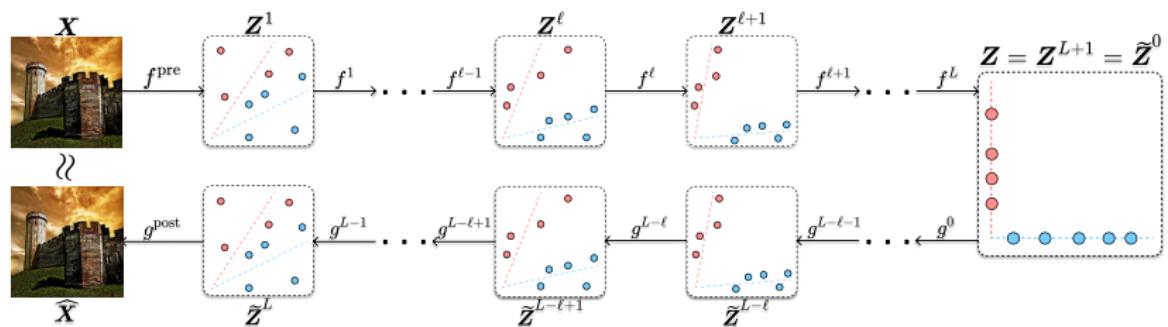
$$\tilde{\mathbf{z}}^{\ell+1} = g^\ell(\tilde{\mathbf{z}}^\ell) = \tilde{\mathbf{z}}^\ell + \tau^\ell \nabla \log q^\ell(\tilde{\mathbf{z}}^\ell)$$

Diffusion models are **iterative denoisers**.

# Structured Denoising-Diffusion

Our alternative: **structured denoising-diffusion**, i.e.,

low-dim. data distribution  $\longleftrightarrow$  structured representations.



## Recall: Low-Dimensional Gaussian Mixture Codebook

We introduce a parametric statistical model for the representations:

Given orthonormal bases  $\mathbf{U}_{[K]} = (\mathbf{U}_k)_{k \in [K]} \in (\mathbb{R}^{d \times p})^K$ , we desire that **each token  $z_i$  of the representation  $Z$  has marginal distribution given by**

$$z_i \stackrel{d}{=} \mathbf{U}_{s_i} \boldsymbol{\alpha}_i, \quad \forall i \in [n]$$

where  $(s_i)_{i \in [n]} \in [K]^n$  are random variables corresponding to the subspace indices, and  $(\boldsymbol{\alpha}_i)_{i \in [n]} \in (\mathbb{R}^p)^n$  are zero-mean standard Gaussian variables.

⇒ both the final representation distribution and incremental transformation operators are structured!

## Connection: Compression and Denoising

**Recall:** In CRATE, the MSSA operator approximates a *gradient descent step* on  $R^c \implies$  compression.

### Theorem (Informal)

Under certain verifiable conditions, if

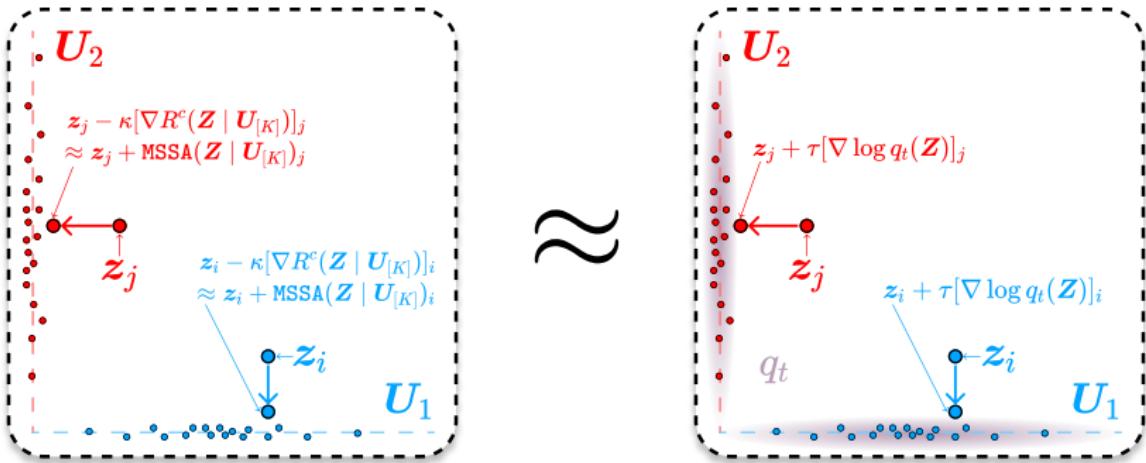
- $\mathbf{Z} \sim$  Gaussian tokens model w/ signal model  $\mathbf{U}_{[K]}$ ,
- and  $\tilde{\mathbf{Z}}_\sigma = \mathbf{Z} + \sigma \mathbf{G}$  where  $\sigma > 0$  is small,  $\mathbf{G} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and  $q^\sigma$  is density of  $\tilde{\mathbf{Z}}_\sigma$ ,

then w.h.p.  $-\nabla R^c(\tilde{\mathbf{Z}}_\sigma \mid \mathbf{U}_{[K]}) \propto \nabla \log q^\sigma(\tilde{\mathbf{Z}}_\sigma)$  (approx.).

### Upshot:

gradient descent on  $R^c \approx$  gradient ascent on score function  
= denoising step

# Compression and Structured Denoising-Diffusion



Structured denoising-diffusion models, where we replace the score with  $-\nabla R^c$ , are **iterative compressors**.

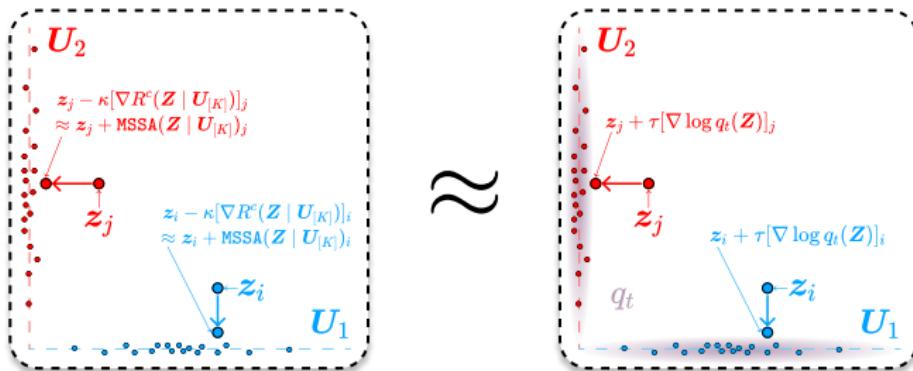
# Structured Denoising-Diffusion

**Structured denoising** process:

$$\frac{d}{dt} \mathbf{Z}(t) = -\frac{1}{2(T-t)} \nabla R^c(\mathbf{Z}(t) | \mathbf{U}_{[K]}(t)), \quad \forall t \in [0, T]$$

Time reversal = **structured diffusion** process:

$$\frac{d}{dt} \tilde{\mathbf{Z}}(t) = \frac{1}{2t} \nabla R^c(\tilde{\mathbf{Z}}(t) | \mathbf{U}_{[K]}(T-t)), \quad \forall t \in [0, T]$$



# Structured Denoising-Diffusion and Network Architectures

Discretizations:

$$\mathbf{Z}^{\ell+} = \mathbf{Z}^\ell - \kappa \nabla R^c(\mathbf{Z}^\ell \mid \mathbf{U}_{[K]}^\ell),$$

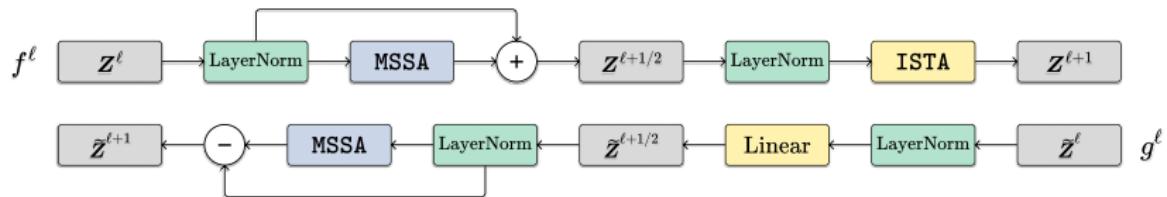
$$\tilde{\mathbf{Z}}^{\ell+} = \mathbf{Z}^\ell + \kappa \nabla R^c(\tilde{\mathbf{Z}}^\ell \mid \mathbf{U}_{[K]}^{L-\ell})$$

Approximations:

$$\mathbf{Z}^{\ell+} \approx (1 - \beta\kappa) \mathbf{Z}^\ell + \beta\kappa \text{MSSA}(\mathbf{Z}^\ell \mid \mathbf{U}_{[K]}^\ell),$$

$$\tilde{\mathbf{Z}}^{\ell+} \approx (1 + \beta\kappa) \tilde{\mathbf{Z}}^\ell - \beta\kappa \text{MSSA}(\tilde{\mathbf{Z}}^\ell \mid \mathbf{U}_{[K]}^{L-\ell})$$

# White-Box Autoencoder



(Encoder)

$$\mathbf{Z}^{\ell+1/2} = \mathbf{Z}^\ell + \text{MSSA}(\mathbf{Z}^\ell \mid \mathbf{U}_{[K]}^\ell)$$

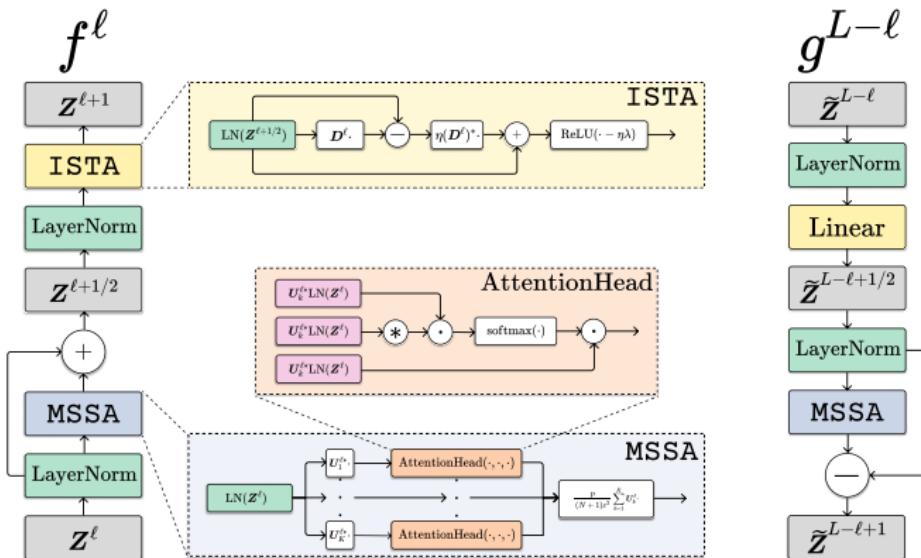
$$\mathbf{Z}^{\ell+1} = \text{ISTA}(\mathbf{Z}^\ell \mid \mathbf{D}^\ell)$$

(Decoder)

$$\tilde{\mathbf{Z}}^{\ell+1/2} = \mathbf{E}^\ell \tilde{\mathbf{Z}}^\ell$$

$$\tilde{\mathbf{Z}}^{\ell+1} = \tilde{\mathbf{Z}}^{\ell+1/2} - \text{MSSA}(\tilde{\mathbf{Z}}^{\ell+1/2} \mid \mathbf{V}_{[K]}^\ell)$$

# CRATE-AE Operational Summary

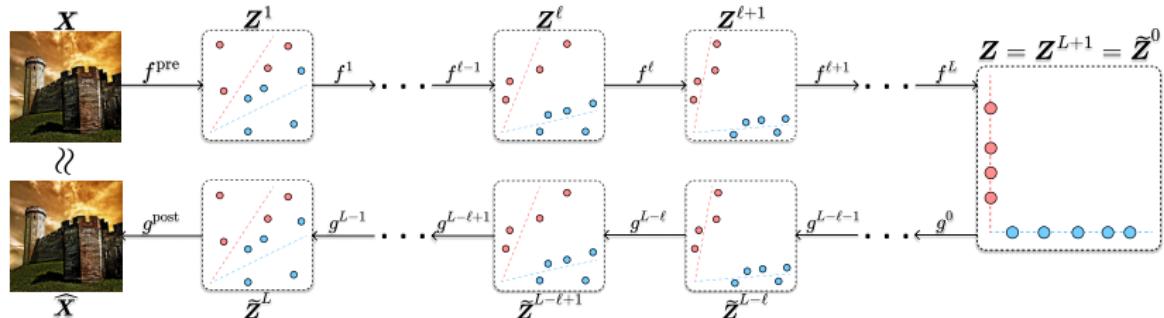


**Encoder:** Same as CRATE.

**Decoder:** MSSA is off by  $\times - 1$ , “MLP” is a linear layer.

# Learning Signal Models From Data

How to learn the true signal model(s)?



- Many ways to train this model — supervised/unsupervised
- Is there a training loss that enforces that the network models the structure of the data distribution?
- Revisit classical tools for learning + verifying low-dimensional structure: denoising, compression, **matrix completion**

# Low-Rank Matrix Completion

**Task:** Given  $\mathbf{X} = \text{Mask}_\Omega(\mathbf{X}_\natural)$ , reconstruct  $\mathbf{X}_\natural$ .

**Method:** If  $\mathbf{X}_\natural$  has low-rank, then we can efficiently recover it by *nuclear norm minimization*:

$$\min_{\mathbf{Y}} \left[ \|\mathbf{Y}\|_* \doteq \sum_i \sigma_i(\mathbf{Y}) \right] \quad \text{subject to} \quad \mathbf{X} = \text{Mask}_\Omega(\mathbf{Y}).$$

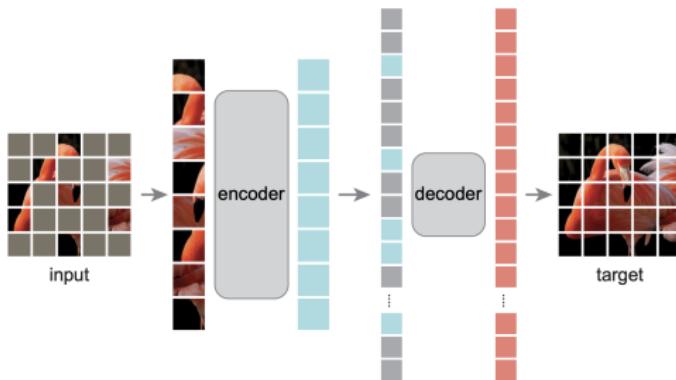
$$\begin{array}{c} \text{Users} \\ \vdots \\ \text{Icon} \end{array} \begin{bmatrix} 5 & 3 & \dots & ? \\ ? & 2 & \dots & 4 \\ \vdots & \vdots & \ddots & \vdots \\ 5 & ? & \dots & ? \end{bmatrix} \quad = \quad \mathcal{P}_\Omega \begin{pmatrix} \begin{bmatrix} 5 & 3 & \dots & 5 \\ 4 & 2 & \dots & 4 \\ \vdots & \vdots & \ddots & \vdots \\ 5 & 5 & \dots & 3 \end{bmatrix} \\ \text{Complete Ratings } \mathbf{X} \end{pmatrix}$$

Items  
Observed (Incomplete) Ratings  $\mathbf{Y}$

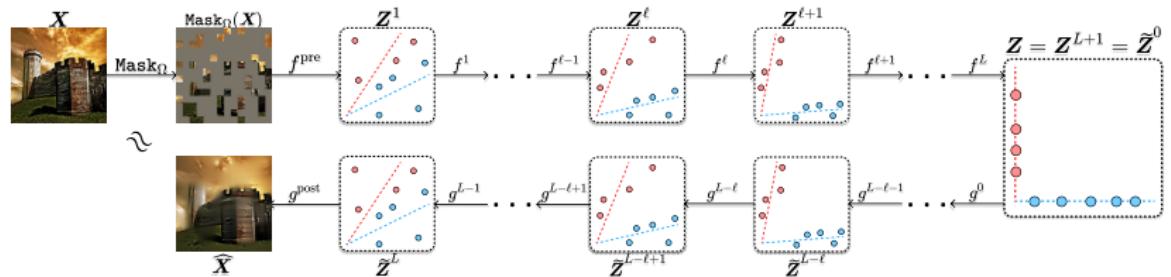
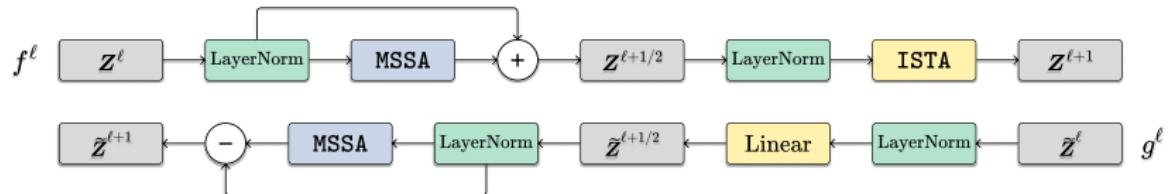
# Learning Low-Dimensional Structure via MAE

**Question:** What if the hidden  $X_h$  has *nonlinear, learnable* low-dimensional structure?

Modern, nonlinear, deep analogue: **masked autoencoding**.

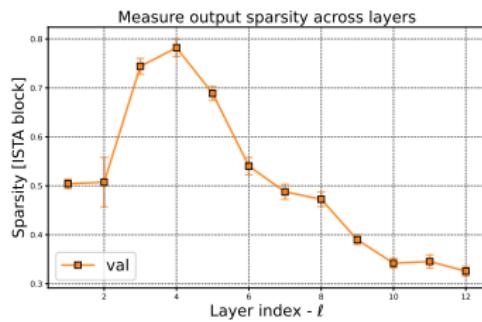
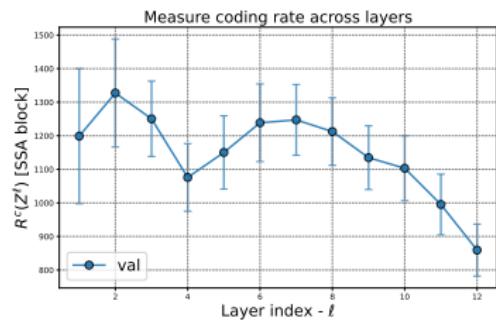


# CRATE-MAE



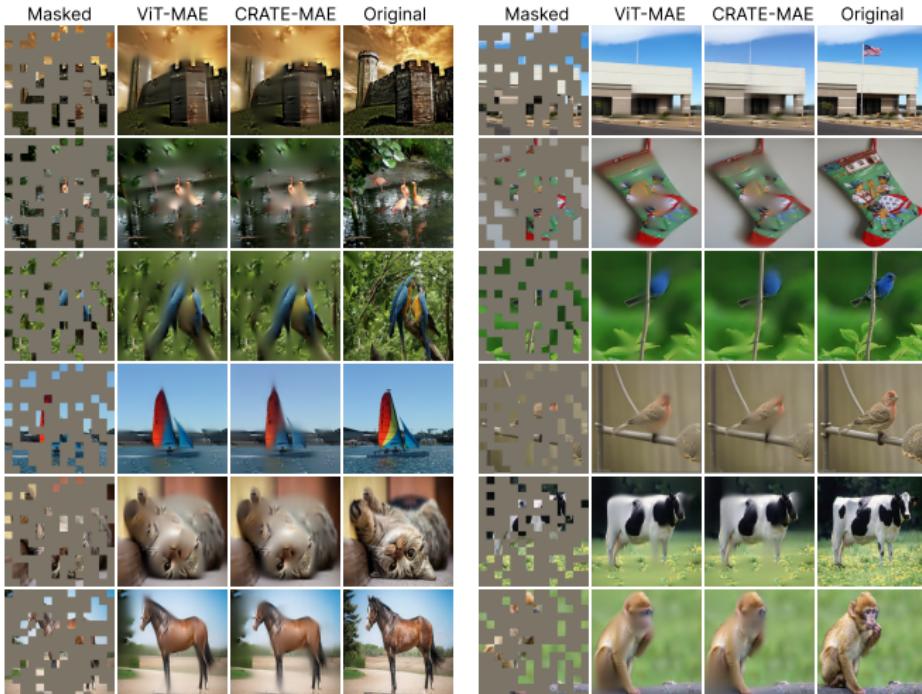
$$\mathcal{L}(f, g) \doteq \mathbb{E} \|X - g(f(\text{Mask}_\Omega(X)))\|_2^2$$

# CRATE-MAE: Experimental Validation



Encoder learns **compressed and sparse** representation.

# CRATE-MAE: Experimental Validation



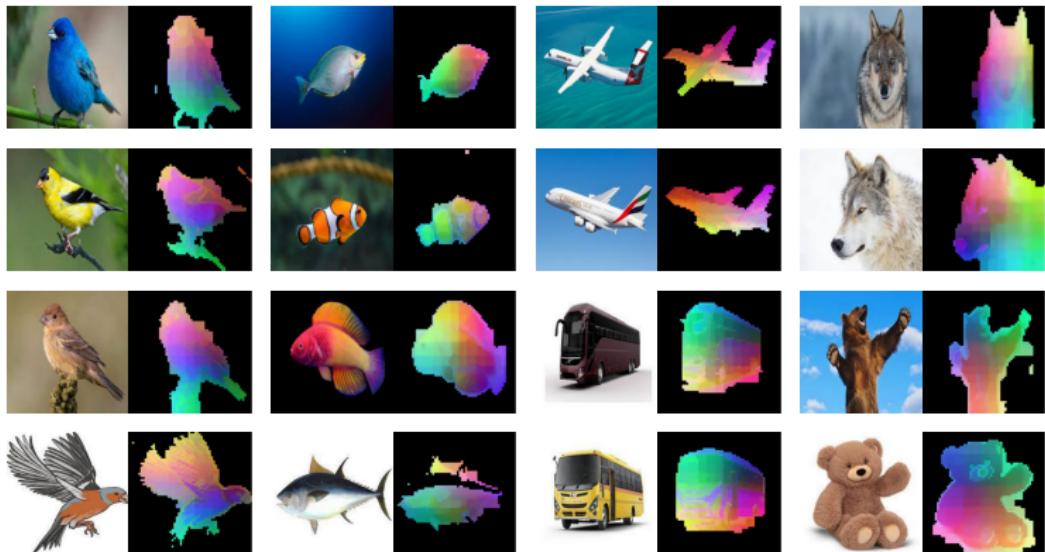
Visually similar reconstructions compared to ViT-MAE.

# CRATE-MAE: Experimental Validation

Model	CRATE-MAE-S	CRATE-MAE-B	ViT-MAE-S	ViT-MAE-B
# parameters	25.4M	44.6M	47.6M	143.8M
<u>Fine-Tuning</u>				
CIFAR10	96.2	96.8	97.6	98.5
CIFAR100	79.0	80.3	83.0	87.0
Oxford Flowers-102	71.7	78.5	84.2	92.5
Oxford-IIIT-Pets	73.7	76.7	81.7	90.3
<u>Linear Probing</u>				
CIFAR10	79.4	80.9	79.9	87.9
CIFAR100	56.6	60.1	62.3	68.0
Oxford Flowers-102	57.7	61.8	66.8	66.4
Oxford-IIIT-Pets	40.6	46.2	51.8	80.1

Similar feature learning performance as ViT-MAE.

# CRATE-MAE: Experimental Validation



Features from the same class have **linear subspace structure**.

# Outline

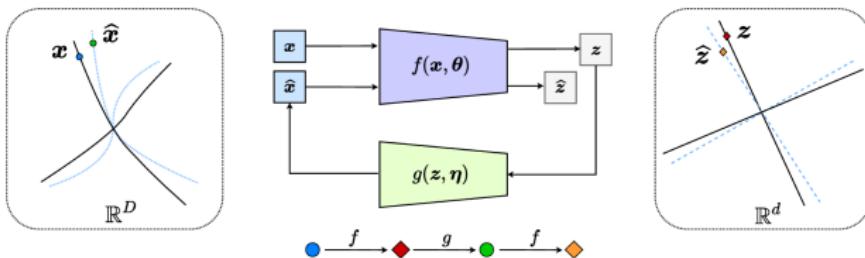
White-Box Autoencoders via Structured Denoising-Diffusion

Closed-Loop Transcription: Self-Consistent Representation Learning

Tutorial Conclusion

# From Open-Loop to Closed-Loop?

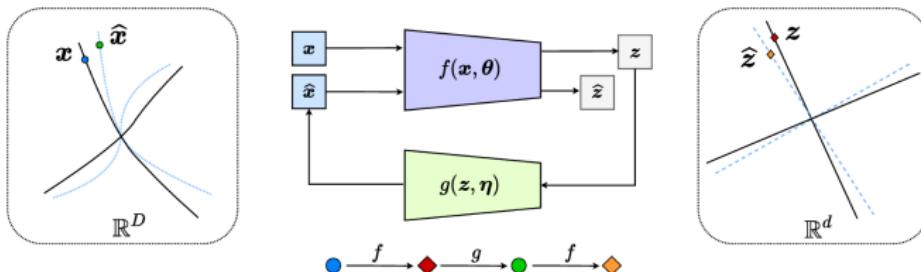
- Start: White-box deep architectures via unrolled optimization
- Recently: Auto-encoding for invertible representation
- Now: Closed-loop learning for self-consistency



How to ensure that the representations/signal models  
are learned *consistently* and *autonomously*?

# From Open-Loop to Closed-Loop

- Open-loop white-box networks: need human-provided learning signal, i.e., class label (for classification task)
- Bi-directional white-box autoencoders: learning signal given by distance in  $\mathcal{X}$  space, i.e., data space
- Closed-loop: learning signal given by distance in  $\mathcal{Z}$  space between representations  $\Rightarrow$  self-consistency



# From Auto-Encoding to Closed-Loop

Autoencoding loss:

$$\mathcal{L}(f, g) \doteq \mathbb{E}[d_{\ell^2}(\mathbf{X}, \widehat{\mathbf{X}})^2], \text{ where } \widehat{\mathbf{X}} = g(f(\mathbf{X})).$$

Masked autoencoding loss:

$$\mathcal{L}(f, g) \doteq \mathbb{E}[d_{\ell^2}(\mathbf{X}, \widehat{\mathbf{X}})^2], \text{ where } \widehat{\mathbf{X}} = g(f(\text{Mask}_{\Omega}(\mathbf{X}))).$$

Wasserstein GAN loss:

$$\mathcal{L}(f, g) \doteq d_{W_1}(\text{Law}(\mathbf{X}), \text{Law}(\widehat{\mathbf{X}}) \mid f), \text{ where } \widehat{\mathbf{X}} = g(\mathbf{W}).$$

**What is the right distance  $d(X, \widehat{X})$  for images?**

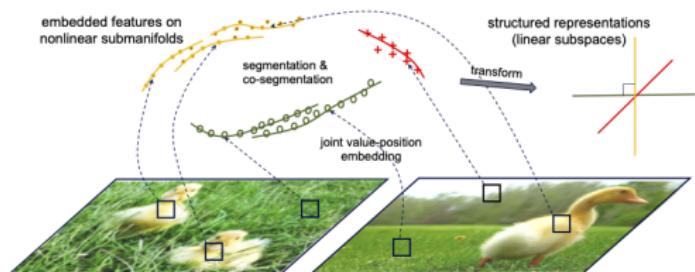
# From Auto-Encoding to Closed-Loop

*Rate reduction* provides a natural measure of difference between representations  $Z_1, Z_2$ :

$$\Delta R(Z_1, Z_2) = R([Z_1, Z_2]) - \frac{1}{2}(R(Z_1) + R(Z_2)).$$

A good measure for representations  $Z$ , but what about images  $X$ ?

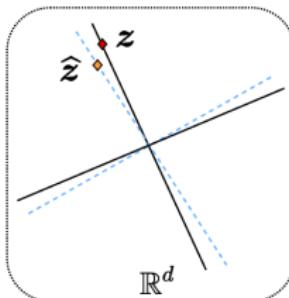
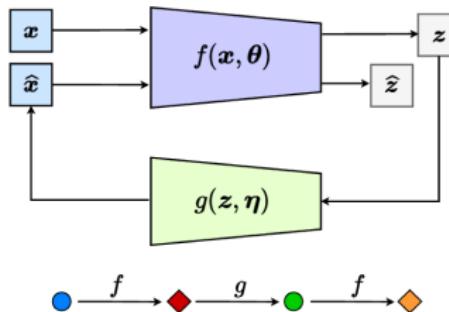
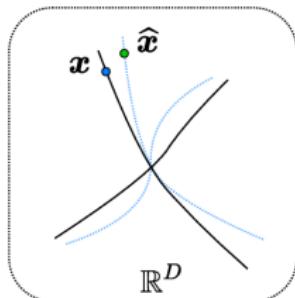
$$X \xrightarrow{f} Z \xrightarrow{g} \widehat{X}$$



Do we ever need to measure distance in  $X$  space?

# From Auto-Encoding to Closed-Loop

Can we measure distances **only** in  $Z$  space? **Yes!**



$$X \xrightarrow{f} Z \xrightarrow{g} \widehat{X} \xrightarrow{f} \widehat{Z}$$

Measure difference between data  $X$  and autoencoding  $\widehat{X}$  through their representations  $Z$  and  $\widehat{Z}$ :

$$\Delta R(Z, \widehat{Z}) = R([Z, \widehat{Z}]) - \frac{1}{2}(R(Z) + R(\widehat{Z})).$$

# Closed-Loop Transcription

This is the **closed-loop transcription (CTRL) framework**.

What are the roles of  $f$  and  $g$ , and how to train a model?

- $f$  acts as an *encoder, discriminator, and sensor*.
- $g$  acts as a *decoder, generator, and controller*.

A *closed-loop* notion of **self consistency** between  $\mathbf{X}$  and  $\widehat{\mathbf{X}}$  is defined by  $f(\mathbf{X}) \approx f(\widehat{\mathbf{X}})$  and achieved by a *minimax game*:

$$\max_f \min_g \Delta R(\mathbf{Z}, \widehat{\mathbf{Z}}) \quad \text{where } \mathbf{Z} = f(\mathbf{X}), \widehat{\mathbf{Z}} = f(g(\mathbf{Z}))$$

## CTRL: Theory for Special Case

### CTRL, Linear Case (Informal)

Suppose columns of  $\mathbf{X}$  are drawn from (possibly correlated) linear subspaces. Then, under mild conditions (e.g., on dimensions), if  $f$  and  $g$  are linear functions, at the CTRL minimax game's equilibrium we have  $\text{span}(\mathbf{Z}) = \text{span}(\widehat{\mathbf{Z}})$ .

**Upshot:** CTRL learns self-consistent representations if data is piecewise linear (where linear functions are all you need).

Expect that this generalizes to more realistic settings (data supported on nonlinear manifolds, statistical assumptions, etc.), motivating its use in practice with NNs.

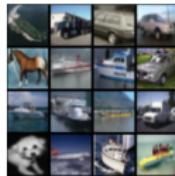
# CTRL: Qualitative Generative Results

5	0	4	1	9	2	1	3
1	4	3	5	3	6	1	7
2	8	6	9	4	0	9	1
1	2	4	3	2	7	3	8
6	9	0	5	6	0	7	4
1	8	1	9	3	9	8	5
3	3	0	7	7	9	8	8
0	9	4	1	4	4	6	0

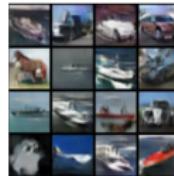
(a) MNIST  $\mathbf{X}$

5	0	4	1	9	2	1	3
1	4	3	5	3	6	1	7
2	8	6	9	4	0	9	1
1	2	4	3	2	7	3	8
6	9	0	5	6	0	7	4
1	8	1	9	3	9	8	5
3	3	0	7	7	9	8	8
0	9	4	1	4	4	6	0

(b) MNIST  $\hat{\mathbf{X}}$



(c) CIFAR-10  $\mathbf{X}$



(d) CIFAR-10  $\hat{\mathbf{X}}$



(e) ImageNet  $\mathbf{X}$



(f) ImageNet  $\hat{\mathbf{X}}$

CTRL maps semantically similar images to similar features.

⇒ **self-consistent representation learning!**

# CTRL: Quantitative Generative Results

Method	CIFAR-10		STL-10		ImageNet	
	IS↑	FID↓	IS↑	FID↓	IS↑	FID↓
<i>GAN based methods</i>						
DCGAN ( <a href="#">Radford et al., 2015</a> )	6.6	-	7.8	-	-	-
SNGAN ( <a href="#">Miyato et al., 2018</a> )	7.4	29.3	<b>9.1</b>	40.1	-	48.73
CSGAN ( <a href="#">Wu et al., 2019b</a> )	8.1	19.6	-	-	-	-
LOGAN ( <a href="#">Wu et al., 2019a</a> )	<b>8.7</b>	<b>17.7</b>	-	-	-	-
<i>VAE/GAN based methods</i>						
VAE ( <a href="#">Kingma &amp; Welling, 2013</a> )	3.8	115.8	-	-	-	-
VAE/GAN ( <a href="#">Larsen et al., 2016</a> )	7.4	39.8	-	-	-	-
NVAE ( <a href="#">Vahdat &amp; Kautz, 2020</a> )	-	50.8	-	-	-	-
DC-VAE ( <a href="#">Parmar et al., 2021</a> )	<b>8.2</b>	<b>17.9</b>	8.1	41.9	-	-
LDR-Binary (ours)	<b>8.1</b>	<b>19.6</b>	8.4	<b>38.6</b>	7.74	<b>46.95</b>
LDR-Multi (ours)	7.1	23.9	7.7	45.7	6.44	55.51

CTRL is competitive as a generative model vs. setups with similar compute/parameter count.

# CTRL: Disentangled Representations

Hat



Hair Color

Glasses

(a) Disentangled attributes as principal components



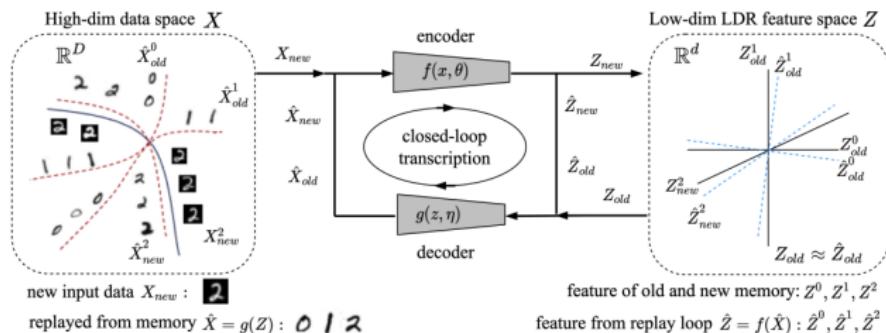
(b) Interpolation between distinct samples

⇒ CTRL learns **disentangled latent representations**.

# Eliminating Catastrophic Forgetting via CTRL

Closed-loop transcription alleviates catastrophic forgetting.

Test case: train a model one class at a time.



$$\begin{aligned} & \max_f \min_g \quad \Delta R(Z) + \Delta R(\hat{Z}) + \Delta R(Z_{\text{new}}, \hat{Z}_{\text{new}}) \\ \text{s.t.} \quad & \Delta R(Z_{\text{old}}, \hat{Z}_{\text{old}}) = 0. \end{aligned}$$

# CTRL: Quantitative Incremental Learning Results

Category	Method	MNIST				CIFAR-10			
		10-splits		5-splits		10-splits		5-splits	
		Last	Avg	Last	Avg	Last	Avg	Last	Avg
<i>Regularization</i>	LwF (Li & Hoiem, 2017)	-	-	0.196	0.455	-	-	0.196	0.440
	SI (Zenke et al., 2017)	-	-	0.193	0.461	-	-	0.196	0.441
<i>Architecture</i>	ReduNet (Wu et al., 2021)	-	-	0.961	0.982	-	-	0.539	0.645
<i>Exemplar</i>	iCaRL (Rebuffi et al., 2017)	0.322	0.588	0.725	0.803	0.212	0.431	0.487	0.632
	A-GEM (Chaudhry et al., 2018)	0.382	0.574	0.597	0.764	0.115	0.293	0.204	0.473
	CLR-ER (Arani et al., 2022)	-	-	-	0.895	-	-	-	0.662
<i>Generative</i>	DGMw (Ostapenko et al., 2019)	-	0.965	-	-	-	0.562	-	-
	EEC (Ayub & Wagner, 2020)	-	0.978	-	-	-	0.669	-	-
	EECS (Ayub & Wagner, 2020)	-	0.963	-	-	-	0.619	-	-
	i-CTRL (ours)	<b>0.975</b>	<b>0.989</b>	<b>0.978</b>	<b>0.990</b>	<b>0.599</b>	<b>0.727</b>	<b>0.627</b>	<b>0.723</b>

CTRL framework enables powerful incremental learning.  
⇒ **autonomous, continual learning?**

## Key References

- Masked Completion via Structured Diffusion with White-Box Transformers  
<https://arxiv.org/abs/2404.02446>
- Closed-Loop Data Transcription to an LDR via Minimaxing Rate Reduction  
<https://arxiv.org/abs/2111.06636>
- Website: <https://ma-lab-berkeley.github.io/CRATE/>
- GitHub: <https://github.com/Ma-Lab-Berkeley/CRATE>

# Outline

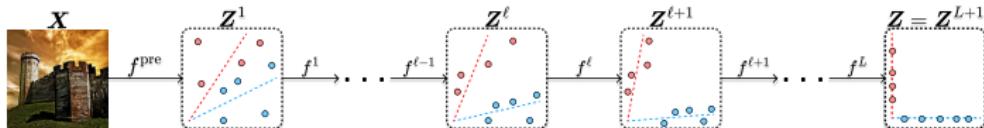
White-Box Autoencoders via Structured Denoising-Diffusion

Closed-Loop Transcription: Self-Consistent Representation Learning

Tutorial Conclusion

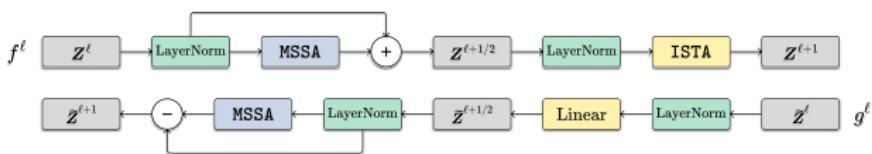
# Conclusions of This Tutorial, Part 1

- **White-box deep networks via rate reduction.**
  - Introduce a principled objective which promotes *compressed*, *linearized* and *sparse* representations.
  - Then, deep networks, such as ReduNet ( $\approx$  ResNet), are just unrolled iterative schemes for optimizing such an objective.
- **White-box transformer via sparse rate reduction.**
  - A principled objective which promotes *compressed*, *linearized*, and *sparse* representations.
  - A white-box transformer architecture, CRATE, constructed from unrolled optimization on the sparse rate reduction, which performs compressive encoding at scale.

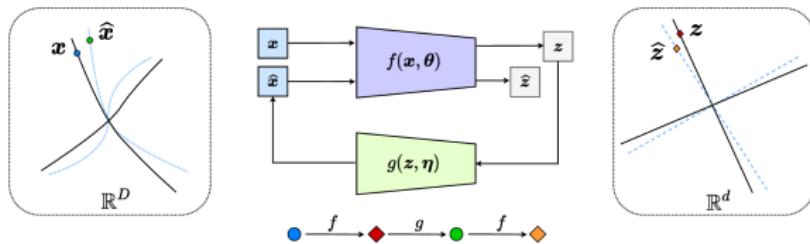


# Conclusions of This Tutorial, Part 2

- **Bi-directional autoencoding for consistency.**
  - Structured denoising-diffusion (unifying compression, denoising, and transformers in the case where the desired representation is a mixture of subspaces).

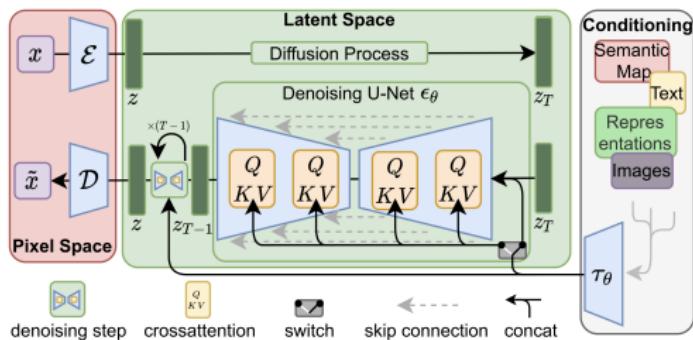


- **Closed-loop transcription for self-consistency.**
  - Autonomous self-consistent representation learning (provably correct for a mixture of subspaces).



# What's Next?

- More efficient architectures using more advanced optimization and sampling techniques.
- Better learning strategies (layer-wise? backward and forward?)
- Truly controllable generation for generative models which also exploit learned compact and structured representations.



# Tutorial Conclusion

---

**Thank You! Questions?**