

# Learning Deep Low-Dimensional Models from High-Dimensional Data: From Theory to Practice

(White-Box Transformers via Sparse Rate Reduction)

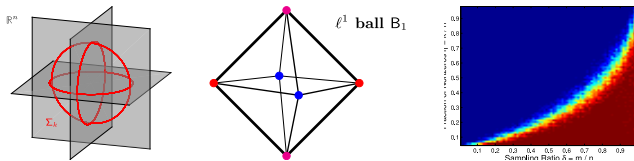
**Sam Buchanan**

University of California, Berkeley

October 19, 2025



# High-Dimensional Data Analysis: Sparse Reconstruction



**Sparse recovery:** **structured** signals, **linear** measurements

$$\mathbf{x} = \mathbf{A}\mathbf{z}_o, \quad \mathbf{z}_o \text{ sparse}, \quad \mathbf{A} \in \mathbb{R}^{m \times n} \text{ random}$$

with **convex** optimization

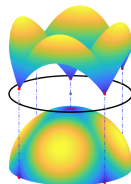
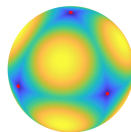
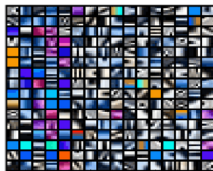
$$\mathbf{z}_\star = \arg \min_{\mathbf{z} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{x} - \mathbf{A}\mathbf{z}\|_2^2 + \lambda \|\mathbf{z}\|_1$$

and provable (high probability) guarantees

$$\mathbf{z}_\star = \mathbf{z}_o \text{ when } \text{measurements} \gtrsim \text{sparsity} \times \log \left( \frac{\text{dimension}}{\text{sparsity}} \right)$$



# Representation (Dictionary) Learning



Dictionary learning: **structured** signals, **bilinear** measurements

$$\mathbf{X} = \mathbf{A}_o \mathbf{Z}_o \in \mathbb{R}^{n \times p}, \quad \mathbf{Z}_o \text{ sparse and random}, \quad \mathbf{A}_o^* \mathbf{A}_o \approx \mathbf{I}$$

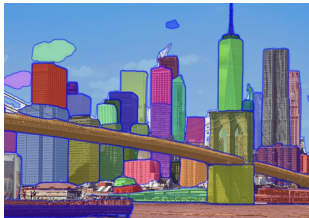
with (efficient) **nonconvex** optimization

$$\mathbf{a}_\star = \arg \min_{\|\mathbf{a}\|_2=1} \|\mathbf{X}^* \mathbf{a}\|_1$$

and provable (high probability) guarantees

$$\mathbf{a}_\star \approx (\mathbf{A}_o)_j \text{ when } \text{observations} \geq \text{poly}(\text{expected sparsity})$$

# Modern (Deep) Representation Learning



Perceiving the physical world  $\Rightarrow$  **nonlinear signals!**  
 Nonlinearity demands **deeper** representations.



**In-the-Wild Data**  
 Over 4.1 million images  
 Five diverse data sources



**Masked Autoencoder**

(a) Masking (b) Autoencoder

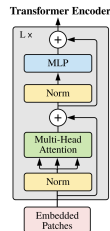
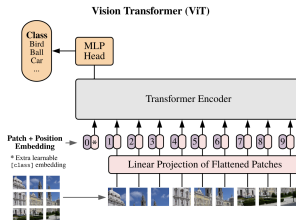
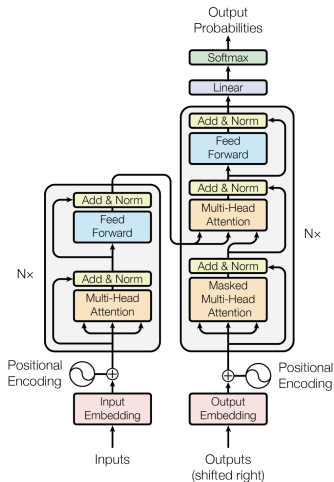


**Real-World Robotic Tasks**

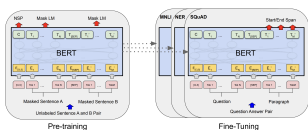
Two robots (xArm, Allegro hand)  
 Eight tasks (scenes, objects)



# Transformers: Modern Representation Learning's Workhorse



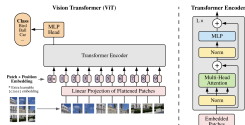
# Transformers: A Universal Backbone



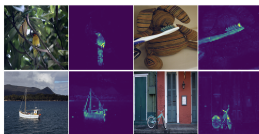
BERT



GPT



ViT



DINO



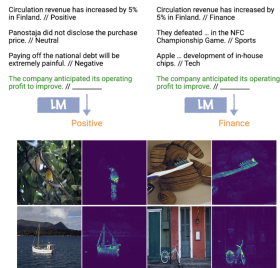
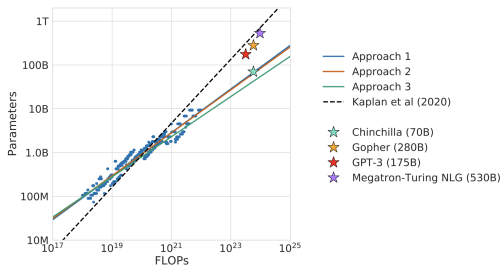
TF + NLP

TF + Vision

TF + Robotics

# Shortcomings of Black-Box Models?

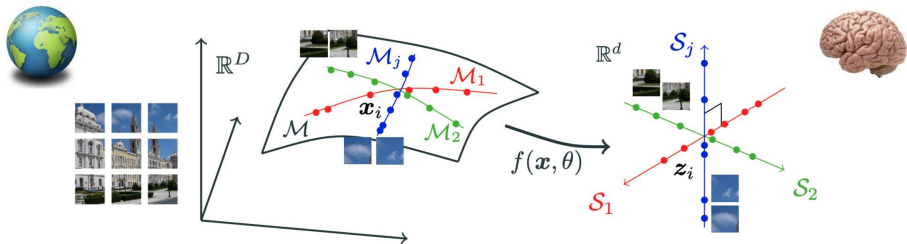
Transformers are **empirically-designed** (or “**black-box**” models).



**How to understand such “emergent” phenomena?**  
**What to do when things go wrong?**

Competing partial theoretical interpretations, e.g. [Vidal 2022] [Bai et al. 2023]  
 [Geshkovski et al. 2023]

# Representations: What and How to Learn?



## The main objective of learning:

Identify **low-dimensional structures** in sensed data of the world and transform to a **compact and structured** representation.

# Outline

## 1 Analytical Models

- Geometry and Sparsity

- Optimization and Neural Networks

## 2 Deep Representation Learning

- Transformers for Visual Data

- Objectives for Representation Learning

- Unrolled Optimization for Representation Learning

- Compression and Self-Attention

- Sparsification and MLP

- Coding Rate Reduction Transformer

- Experimental Results on CRATE

## 3 Conclusions for the Tutorial

# Outline

## 1 Analytical Models

Geometry and Sparsity

Optimization and Neural Networks

## 2 Deep Representation Learning

Transformers for Visual Data

Objectives for Representation Learning

Unrolled Optimization for Representation Learning

Compression and Self-Attention

Sparsification and MLP

Coding Rate Reduction Transformer

Experimental Results on CRATE

## 3 Conclusions for the Tutorial



# A Low-Dimensional Subspace

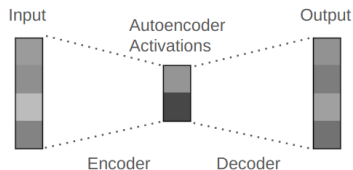
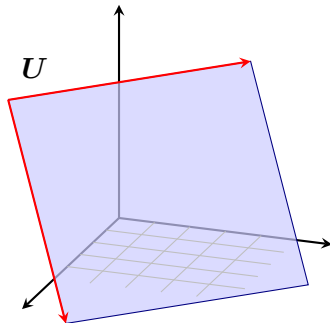
Canonical example:

subspace  $U \in \mathbb{R}^{D \times d}$ :

$$\mathbf{x} \xrightarrow{f=U^\top} \mathbf{z} \xrightarrow{g=U} \hat{\mathbf{x}}$$

Principal component analysis:

$$\min_U \mathbb{E} \left[ \left\| \mathbf{x} - UU^\top \mathbf{x} \right\|_2^2 \right]$$



# Sparsity and Sparse Coding

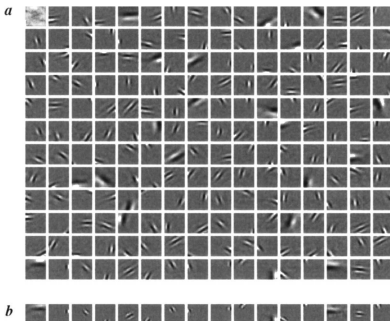
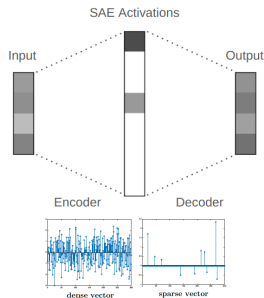
A significant generalization: unions of subspaces

$\ell^0$  “norm”: **# of nonzero entries**,  $\|x\|_0 = |\{i \mid x_i \neq 0\}|$ .

Given (learned)  $A \in \mathbb{R}^{D \times d}$ , represent  $x$  as a sparse code:

$$f(x) = \min_z \|x - Az\|_2^2 + \|z\|_0;$$

$$x \xrightarrow{f} z \xrightarrow{g=A} \hat{x}$$



# How to Learn: Optimization for Low-Dim Structures

We can compute sparse coding with (proximal) gradient descent

$$f(\mathbf{x}) = \arg \min_{\mathbf{z} \geq 0} \|\mathbf{x} - \mathbf{A}\mathbf{z}\|_2^2 + \lambda \|\mathbf{z}\|_1$$

- 1 Given the current code  $\mathbf{z}^\ell$ , gradient descent to better fit  $\mathbf{x}$ ;
- 2 Without moving too much, sparsify the updated code

Then  $f(\mathbf{x}) = \mathbf{z}^\infty$ , where

$$\mathbf{z}^{\ell+1} = \text{ReLU} \left( \eta \mathbf{A}^\top \mathbf{x} + \left( \mathbf{I} - \eta \mathbf{A}^\top \mathbf{A} \right) \mathbf{z}^\ell - \lambda \eta \mathbf{1} \right)$$

# Unrolled Optimization: From Objectives to Deep Networks

Recall the sparse coding objective:

$$f(\mathbf{x}) = \arg \min_{\mathbf{z} \geq 0} \|\mathbf{x} - \mathbf{A}\mathbf{z}\|_2^2 + \lambda \|\mathbf{z}\|_1$$

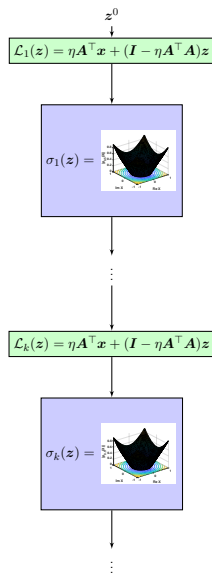
Then  $f(\mathbf{x}) = \mathbf{z}^\infty$ , where

$$\mathbf{z}^{\ell+1} = \text{ReLU} \left( \eta \mathbf{A}^\top \mathbf{x} + \left( \mathbf{I} - \eta \mathbf{A}^\top \mathbf{A} \right) \mathbf{z}^\ell - \lambda \eta \mathbf{1} \right)$$

**Truncate** the network, and **learn** its parameters using data. ( $\mathbf{A} \rightarrow \mathbf{A}^\ell$ )

This approach is called LISTA [Gregor and Lecun, 2010].

$\Rightarrow$  each layer learns its own dictionary!



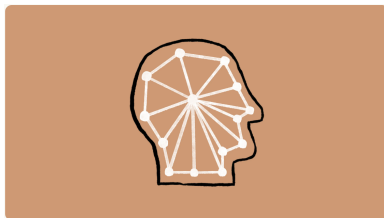
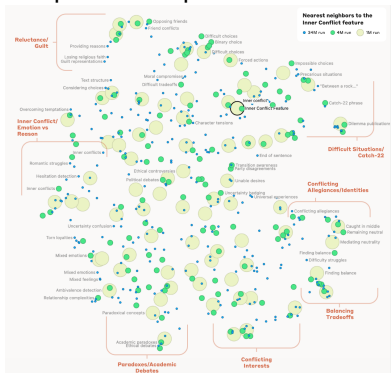
## Special Case: Sparse Autoencoders

Truncate after one iteration & learn a dictionary  $D$  for decoding:

$$f(\mathbf{x}) = \text{ReLU} \left( \mathbf{A}^\top \mathbf{x} - \mathbf{b} \right)$$

$$x \xrightarrow{f} z \xrightarrow{g=D} \hat{x}$$

## Interpretable representations from massive-scale models!



1M/1013764 Code error

[illegible]

Templeton, Conerly et al. 2024

# Using Unrolled Optimization for Deep Learning

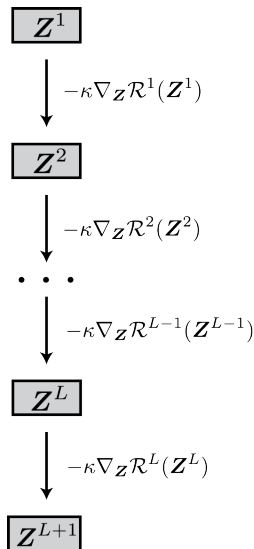
## Unrolled optimization:

- Given objective function  $\mathcal{R}^\ell$ , improve it on input  $\mathbf{Z}^\ell$  by taking optimization step:

$$\mathbf{Z}^{\ell+1} \leftarrow \mathbf{Z}^\ell - \kappa \nabla_{\mathbf{Z}} \mathcal{R}^\ell(\mathbf{Z}^\ell)$$

(...or similar)

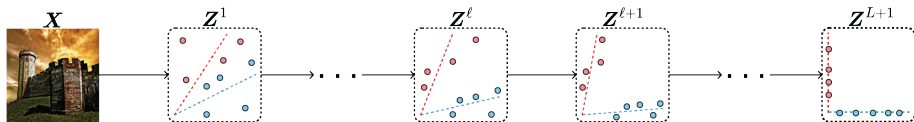
- Collection of objective functions  $(\mathcal{R}^\ell)_{\ell=1}^L$  + optimization strategies  
 $\implies$  data processing algorithm
- New: collection of objective functions + optimization strategies  
 $\implies$  deep network architecture!



# From Unrolled Optimization to Deep Architectures

## Constructing deep networks:

Design objectives  $\mathcal{R}^\ell$  and optimization strategies s.t. unrolling yields  
*compact & structured* deep representation!



**Previously:** Did this for ResNets (ReduNet).

**Next:** *How do we do this for transformers?* What does it buy us?

# Outline

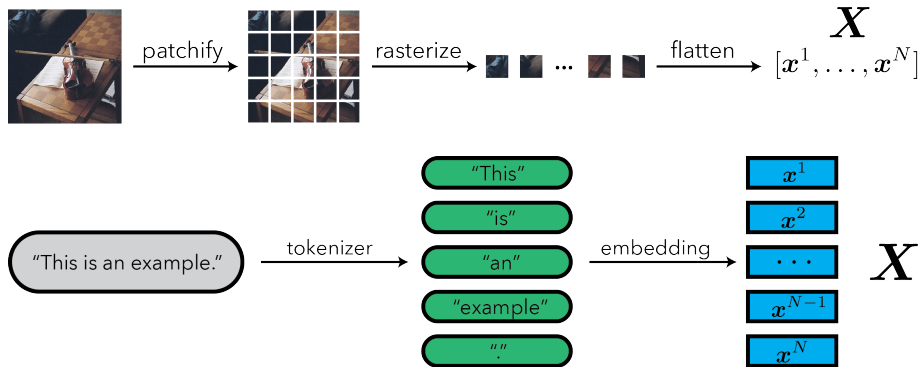
- 1 Analytical Models
  - Geometry and Sparsity
  - Optimization and Neural Networks
- 2 Deep Representation Learning
  - Transformers for Visual Data
  - Objectives for Representation Learning
  - Unrolled Optimization for Representation Learning
  - Compression and Self-Attention
  - Sparsification and MLP
  - Coding Rate Reduction Transformer
  - Experimental Results on CRATE
- 3 Conclusions for the Tutorial



# Scaling Data Processing

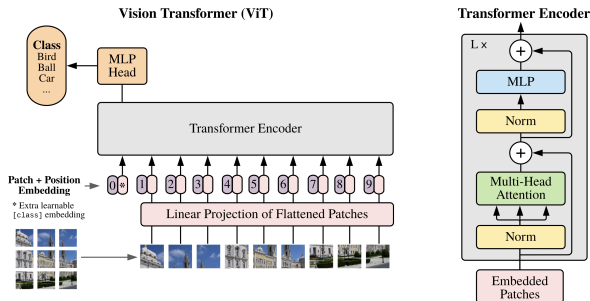
## Data format:

Sequences of *tokens*  $\rightarrow$  *embeddings*  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$



# Processing Images as Token Sequences with ViT

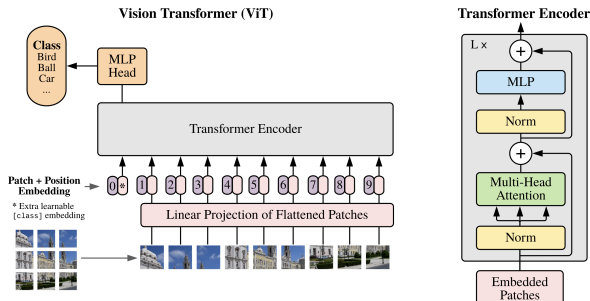
**Recall:** the Vision Transformer (ViT) processes images as a **sequence of patches**.



$$f_{\text{ViT}} = f^L \circ \underbrace{f^{L-1} \circ \dots \circ f^1}_{\text{transformer layer}} \circ \underbrace{f^{\text{pre}}}_{\text{tokenization}}$$

# Processing Images with ViT (Tokenization)

**Recall:** the Vision Transformer (ViT) processes images as a **sequence of patches**.



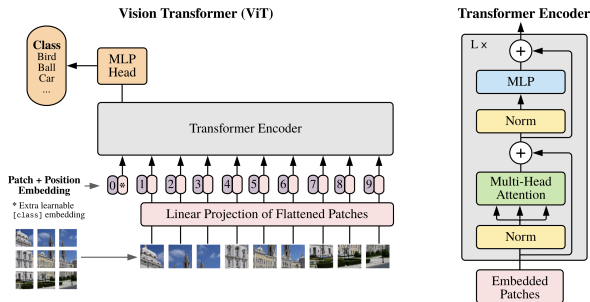
$f^{\text{pre}}$   
tokenization

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$$

$$\mathbf{Z}^1 = [\mathbf{z}_{\text{cls}}, \mathbf{W}^{\text{pre}} \mathbf{X}] + \mathbf{E}_{\text{pos}} = [\mathbf{z}_{\text{cls}}, \mathbf{z}_1, \dots, \mathbf{z}_N] \in \mathbb{R}^{d \times (N+1)}$$

# Processing Images with ViT (TF Block)

**Recall:** the Vision Transformer (ViT) processes images as a **sequence of patches**.



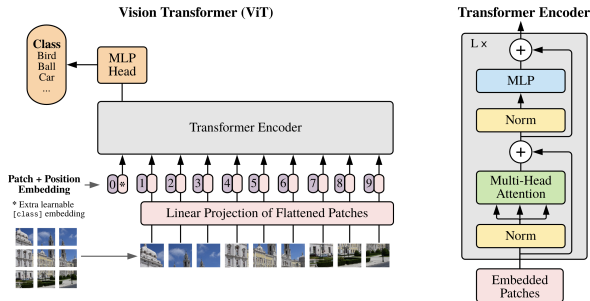
$f^\ell$   
TF layer

$$\mathbf{Z}^{\ell+1/2} = \text{MHSA}(\text{LN}(\mathbf{Z}^\ell)) + \mathbf{Z}^\ell$$

$$f^\ell(\mathbf{Z}^\ell) = \text{MLP}(\text{LN}(\mathbf{Z}^{\ell+1/2})) + \mathbf{Z}^{\ell+1/2}$$

# Processing Images with ViT (TF Block)

**Recall:** the Vision Transformer (ViT) processes images as a **sequence of patches**.



$f^\ell$   
TF layer

$$\mathbf{Z}^{\ell+1/2} = \text{MHSA}(\text{LN}(\mathbf{Z}^\ell)) + \mathbf{Z}^\ell$$

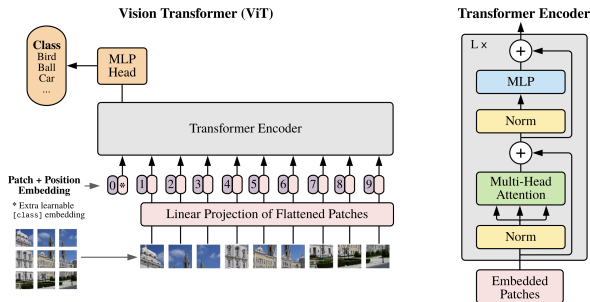
$$f^\ell(\mathbf{Z}^\ell) = \text{MLP}(\text{LN}(\mathbf{Z}^{\ell+1/2})) + \mathbf{Z}^{\ell+1/2}$$

$$\text{SA}(\mathbf{Z}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V) = \mathbf{W}_V \mathbf{Z} \text{softmax}((\mathbf{W}_Q \mathbf{Z})^* (\mathbf{W}_K \mathbf{Z}))$$

$$\text{MHSA}(\mathbf{Z}) = \sum_{h=1}^H \mathbf{W}_{O,h} \text{SA}(\mathbf{Z}; \mathbf{W}_{Q,h}, \mathbf{W}_{K,h}, \mathbf{W}_{V,h})$$

# Processing Images with ViT (TF Block)

**Recall:** the Vision Transformer (ViT) processes images as a **sequence of patches**.



$f^\ell$   
TF layer

$$\mathbf{Z}^{\ell+1/2} = \text{MHSA}(\text{LN}(\mathbf{Z}^\ell)) + \mathbf{Z}^\ell$$

$$f^\ell(\mathbf{Z}^\ell) = \text{MLP}(\text{LN}(\mathbf{Z}^{\ell+1/2})) + \mathbf{Z}^{\ell+1/2}$$

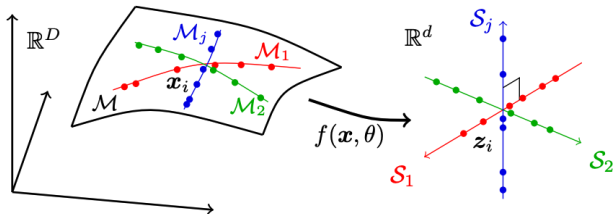
$$\text{MLP}(\mathbf{Z}) = \mathbf{W}_{\text{down}} \sigma(\mathbf{W}_{\text{up}} \mathbf{Z})$$

# Recall (Yi's Lecture): Coding Rate Reduction

**Rate reduction (for non-tokenized data):**

$$\Delta R(\mathbf{Z} \mid \Pi) := R(\mathbf{Z}) - \underbrace{\sum_{k=1}^K \frac{n_k}{n} R(\mathbf{Z}_k)}_{R_c(\mathbf{Z} \mid \Pi)}.$$

Promotes compression of *features (of samples)* against class-wise (learned) low-rank GMM.

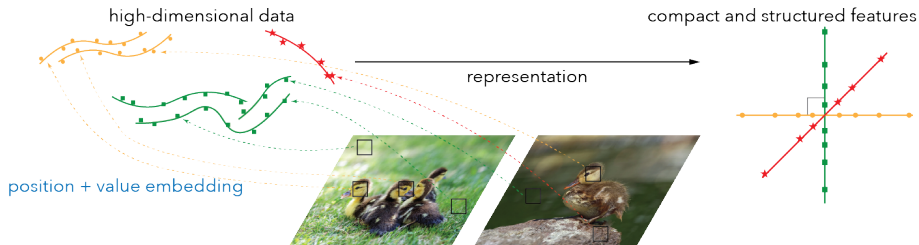


# Rate Reduction for Token Sequences

## Rate reduction for tokenized data:

Parameterize the GMM covariances  $\Sigma_k = \mathbf{U}_k \mathbf{U}_k^\top$ .

$$\Delta R(\mathbf{Z} \mid \underbrace{\mathbf{U}_{[K]}}_{:= (\mathbf{U}_k)_{k=1}^K}) := R(\mathbf{Z}) - \underbrace{\sum_{k=1}^K R(\mathbf{U}_k^\top \mathbf{Z})}_{:= R_c(\mathbf{Z} \mid \mathbf{U}_{[K]})}$$



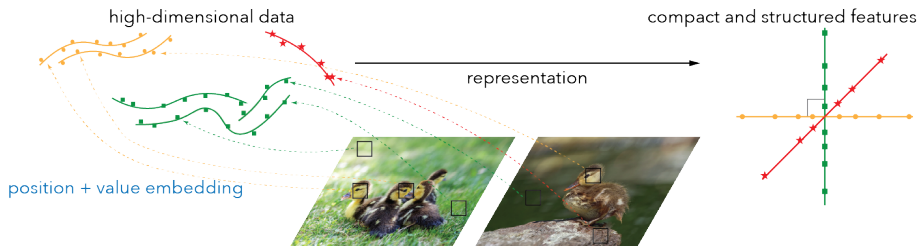


# Sparse Rate Reduction

To be maximally structured, ask  $\mathbf{Z}$  (hence  $\mathbf{U}_k$ ) to be *sparse*!

Objective to maximize: **Sparse Rate Reduction**

$$\text{SRR}(\mathbf{Z} \mid \mathbf{U}_{[K]}) := R(\mathbf{Z}) - R_c(\mathbf{Z} \mid \mathbf{U}_{[K]}) - \lambda \|\mathbf{Z}\|_1$$



# Unrolling the Sparse Rate Reduction

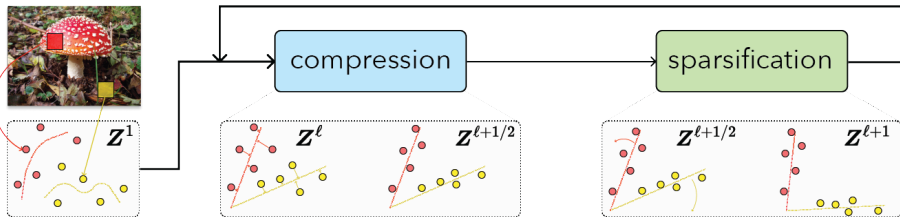
## Proposed optimization strategy:

Two-step (prox-like) iteration.

$$\mathbf{Z}^\ell \mapsto \mathbf{Z}^{\ell+1/2} \mapsto \mathbf{Z}^{\ell+1}$$

$$\mathbf{Z}^{\ell+1/2} \approx \mathbf{Z}^\ell - \kappa \nabla_{\mathbf{Z}} R_c(\mathbf{Z}^\ell \mid \mathbf{U}_{[K]}^\ell) \quad (\text{compression})$$

$$\mathbf{Z}^{\ell+1} \approx \arg \max_{\mathbf{Z}: \mathbf{Z}^{\ell+1/2} = \mathbf{D}^\ell \mathbf{Z}} \{R(\mathbf{Z}) - \lambda \|\mathbf{Z}\|_1\} \quad (\text{sparsification})$$



Parameters:  $(\mathbf{U}_k^\ell)_{k=1}^K \subseteq \mathbb{R}^{d \times p}$ ,  $\mathbf{D}^\ell \in \mathbb{R}^{d \times d}$ .

# Gradient of Compression Objective

Define  $\alpha := p/(n\varepsilon^2)$ .

If  $(U_k)_{k=1}^K \approx \text{orthogonal} + p/w \approx \text{orthogonal} + \approx \text{support } Z$ :

$$\begin{aligned}
 \nabla_Z R_c(Z \mid U_{[K]}) &= \sum_{k=1}^K \alpha (U_k U_k^\top Z) (I_n + \alpha (U_k^\top Z)^\top (U_k^\top Z))^{-1} \\
 &\approx \sum_{k=1}^K \alpha U_k (U_k^\top Z) (I_d - \alpha (U_k^\top Z)^\top (U_k^\top Z)) \\
 &= \alpha \left[ \left( \sum_{k=1}^K U_k U_k^\top \right) Z - \alpha \sum_{k=1}^K U_k (U_k^\top Z) (U_k^\top Z)^\top (U_k^\top Z) \right] \\
 &\approx \alpha \left[ Z - \alpha \sum_{k=1}^K U_k (U_k^\top Z) (U_k^\top Z)^\top (U_k^\top Z) \right]
 \end{aligned}$$

Gradient shaping/"non-parametric autoregression":

$$\nabla_Z R_c(Z) \approx \alpha \left[ Z - \alpha \sum_{k=1}^K U_k (U_k^\top Z) \text{softmax} \left\{ (U_k^\top Z)^\top (U_k^\top Z) \right\} \right]$$

# Multi-head Subspace Self-Attention

$$\nabla_{\mathbf{Z}} R_c(\mathbf{Z}) \approx \alpha \left[ \mathbf{Z} - \alpha \sum_{k=1}^K \mathbf{U}_k (\mathbf{U}_k^\top \mathbf{Z}) \operatorname{softmax} \left\{ (\mathbf{U}_k^\top \mathbf{Z})^\top (\mathbf{U}_k^\top \mathbf{Z}) \right\} \right]$$

**Multi-head Subspace Self-Attention (MSSA):**

$$\text{MSSA}(\mathbf{Z} \mid \mathbf{U}_{[K]}) := \alpha \left[ \mathbf{U}_1, \dots, \mathbf{U}_K \right] \begin{bmatrix} (\mathbf{U}_1^\top \mathbf{Z}) \operatorname{softmax}\{(\mathbf{U}_1^\top \mathbf{Z})^\top (\mathbf{U}_1^\top \mathbf{Z})\} \\ \vdots \\ (\mathbf{U}_K^\top \mathbf{Z}) \operatorname{softmax}\{(\mathbf{U}_K^\top \mathbf{Z})^\top (\mathbf{U}_K^\top \mathbf{Z})\} \end{bmatrix}$$

$$\mathbf{Z}^{\ell+1/2} := \underbrace{(1 - \alpha\kappa) \mathbf{Z}^\ell}_{\text{residual}} + \underbrace{\alpha\kappa \text{MSSA}(\mathbf{Z}^\ell \mid \mathbf{U}_{[K]}^\ell)}_{\text{attention-like}}$$

# Iterative Shrinkage-Thresholding Block

If  $D^\ell \approx$  complete incoherent dictionary then

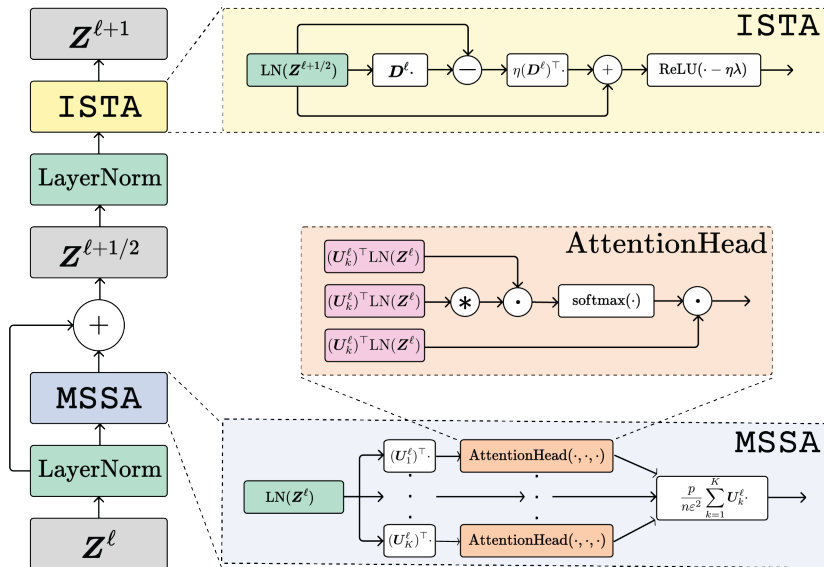
$$\mathbf{Z}^{\ell+1/2} \approx D^\ell \mathbf{Z} \implies R(\mathbf{Z}) \approx R(\mathbf{Z}^{\ell+1/2})$$

Can simplify the prox-like step:

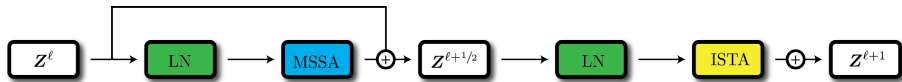
$$\begin{aligned} \mathbf{Z}^{\ell+1} &\approx \arg \max_{\mathbf{Z}: \mathbf{Z}^{\ell+1/2} \approx D^\ell \mathbf{Z}} \{R(\mathbf{Z}) - \lambda \|\mathbf{Z}\|_1\} \approx \arg \min_{\substack{\mathbf{Z} \\ \mathbf{Z}^{\ell+1/2} \approx D^\ell \mathbf{Z}}} \|\mathbf{Z}\|_1 \\ &\approx \arg \min_{\mathbf{Z}} \left\{ \frac{1}{2} \|\mathbf{Z}^{\ell+1/2} - D^\ell \mathbf{Z}\|_2^2 + \lambda' \|\mathbf{Z}\|_1 \right\} \end{aligned}$$

$$\mathbf{Z}^{\ell+1} := \text{ISTA}(\mathbf{Z}^{\ell+1/2}) := \text{ProxGD}(\underbrace{\mathbf{Z}^{\ell+1/2}}_{\text{iterate}}, \underbrace{\mathbf{Z}^{\ell+1/2}}_{\text{target}}, \underbrace{D^\ell}_{\text{dict.}})$$

## CRATE Architecture

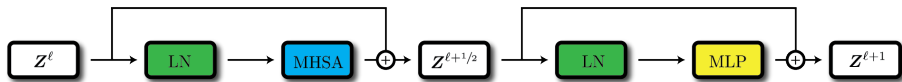


# Comparing CRATE and Regular Transformer



## Three practical differences:

- MSSA sets  $W_{Q,k} = W_{K,k} = W_{V,k} = U_k^\top$
- ISTA sets  $W_{\text{up}} = W_{\text{down}}^\top = D$
- In ISTA the residual connection is moved inside ReLU



# Do CRATE Models Behave According to Theory?

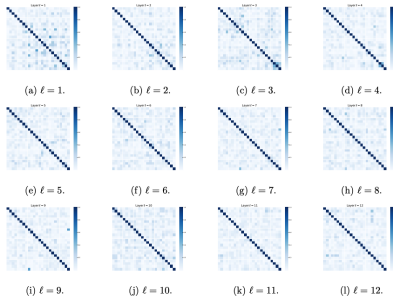
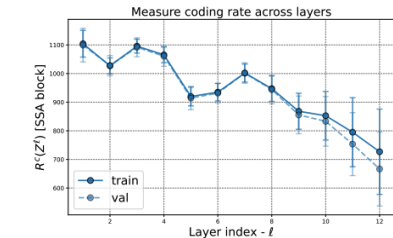


Figure 17: We visualize the  $[U_1^l, \dots, U_K^l]^T [U_1^l, \dots, U_K^l]^T \in \mathbb{R}^{p \times K \times p \times K}$  at different layers. The  $(i, j)$ -th

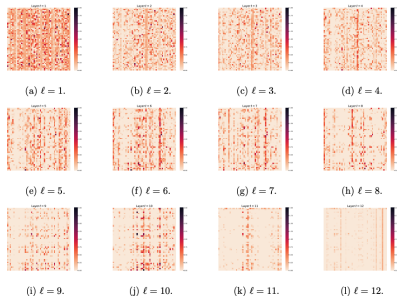
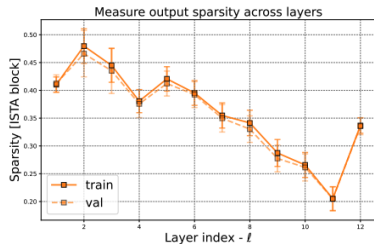


Figure 16: Visualizing layer-wise token  $Z^\ell$  representations at each layer  $\ell$ . To enhance the visual



# Can CRATE Perform Well in Practice?

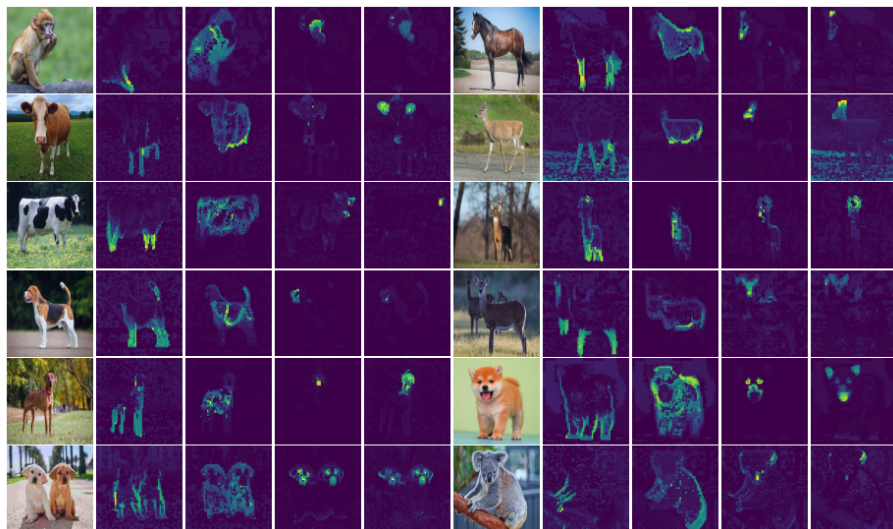
## Vision:

Model	CRATE-T	CRATE-S	CRATE-B	CRATE-L	ViT-T	ViT-S
# parameters	6.09M	13.12M	22.80M	77.64M	5.72M	22.05M
ImageNet-1K	66.7	69.2	70.8	71.3	71.5	72.4
ImageNet-1K ReaL	74.0	76.0	76.5	77.4	78.3	78.4
CIFAR10	95.5	96.0	96.8	97.2	96.6	97.2
CIFAR100	78.9	81.0	82.7	83.6	81.8	83.2
Oxford Flowers-102	84.6	87.1	88.7	88.3	85.1	88.5
Oxford-IIIT-Pets	81.4	84.9	85.3	87.4	88.5	88.6

## Text:

	#parameters	OWT	LAMBADA	WikiText	PTB	Avg
GPT2-Base	124M	2.85	4.12	3.89	4.63	3.87
GPT2-Small	64M	3.04	4.49	4.31	5.15	4.25
CRATE-GPT2-Base	60M	3.37	4.91	4.61	5.53	4.61

# Interpretability and Emergent Segmentation



Head 0  
"Leg"

Head 1  
"Body"

Head 3  
"Face"

Head 4  
"Ear"

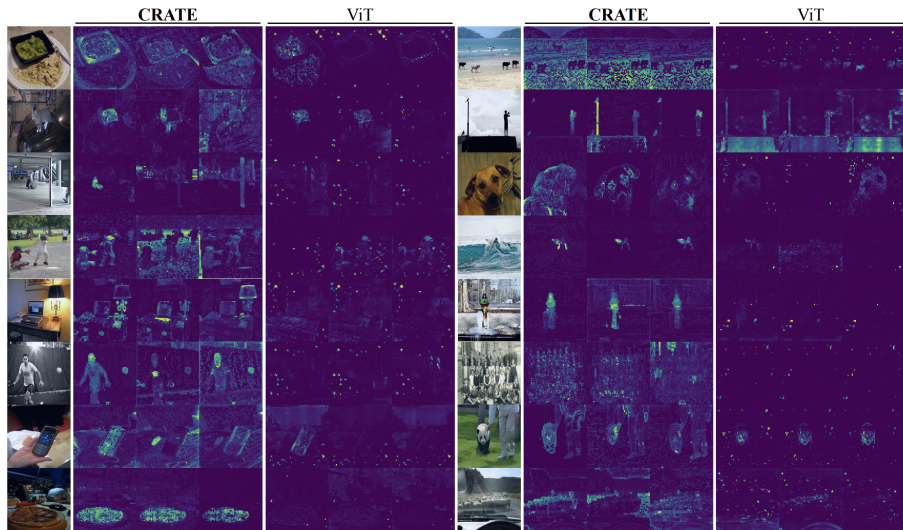
Head 0  
"Leg"

Head 1  
"Body"

Head 3  
"Face"

Head 4  
"Ear"

# Interpretability and Emergent Segmentation



# Performance: Semantic Segmentation

**Setup:** Zero-shot semantic segmentation with CLIP + MSSA.

Top left: original image. Bottom left: CLIP-ViT features. Right: CRATE-CLIP features.



$\approx 5\%$  **better mIoU score than previous approaches!**

# Design Choices in CRATE

## More effective way to do sparsification?

What if we use multiple prox iterations with *overcomplete (wide) dictionary*

$$\mathbf{D}^\ell \in \mathbb{R}^{d \times m}, m > d?$$

$\implies$  different architecture!

### CRATE- $\alpha$ Sparsification Block

$$\begin{aligned}\mathbf{Z}^{\ell+1} &:= \text{ODL}(\mathbf{Z}^{\ell+1/2}) \\ &= \mathbf{D}^\ell \text{ProxGD}(\text{ProxGD}(\mathbf{0}, \mathbf{Z}^{\ell+1/2}, \mathbf{D}^\ell), \mathbf{Z}^{\ell+1/2}, \mathbf{D}^\ell)\end{aligned}$$

# Performance of CRATE- $\alpha$

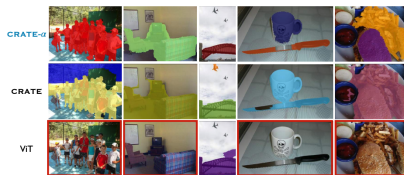
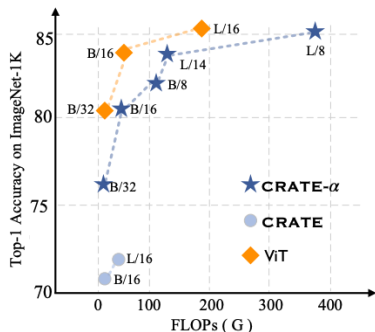
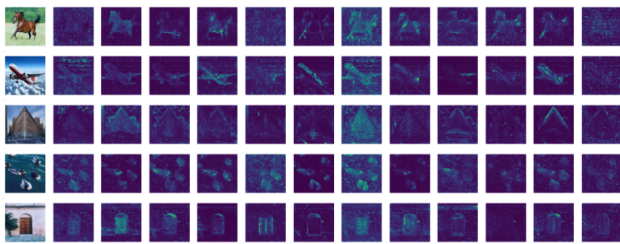


Figure 5: Visualization of segmentation on COCO val2017 [20] with MaskCut [43]. (Top row)

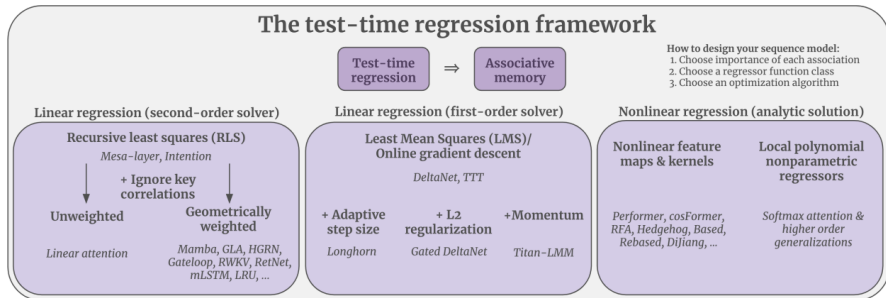
Table 4: The comparison between CRATE and CRATE- $\alpha$  on the NLP task using the OpenWebText dataset.

	GPT-2-base	CRATE-base	CRATE- $\alpha$ -small	CRATE- $\alpha$ -base
Model size	124M	60M	57M	120M
CE val loss	2.85	3.37	3.28	3.14



# Aside: Network Operators as Optimization Primitives

- Optimization gives blocks *similar to* blocks in transformer
- Recent work derives linear attention + similar operators as *exact* optimization steps on regression objectives w.r.t.  $Q, K, V$



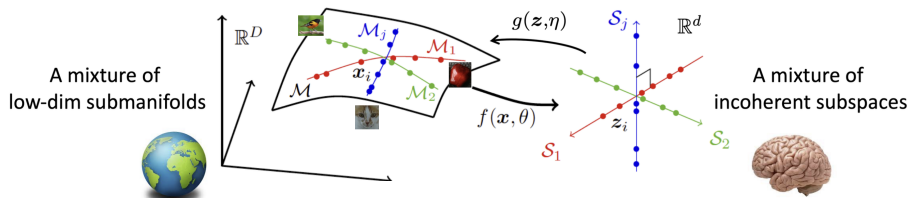
# Outline

- 1 Analytical Models
  - Geometry and Sparsity
  - Optimization and Neural Networks
- 2 Deep Representation Learning
  - Transformers for Visual Data
  - Objectives for Representation Learning
  - Unrolled Optimization for Representation Learning
  - Compression and Self-Attention
  - Sparsification and MLP
  - Coding Rate Reduction Transformer
  - Experimental Results on CRATE
- 3 Conclusions for the Tutorial



# Take-Home Message: Low-Dim Structures are Ubiquitous!

In this tutorial, we have emphasized:

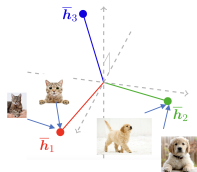
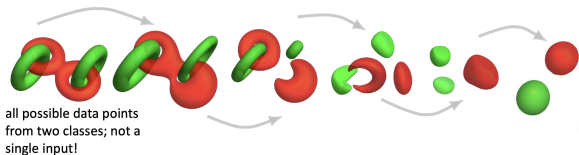


## The objective of learning:

Identify low-dim. distributions in sensed data of the world and transform to a **compact and structured** representation.

**All deep networks** are simply a means to an end!

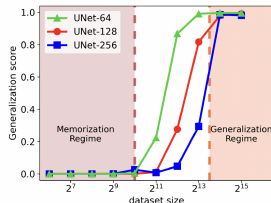
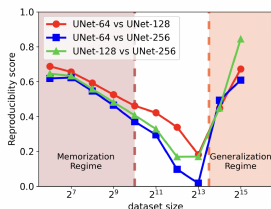
## S2: Understanding Low-Dimensional Structures in Representation Learning



Given enough data and ability to optimize: *inevitable emergence* of low-dim structures in trained deep networks!

Implications for parameter efficiency, transfer learning, ...

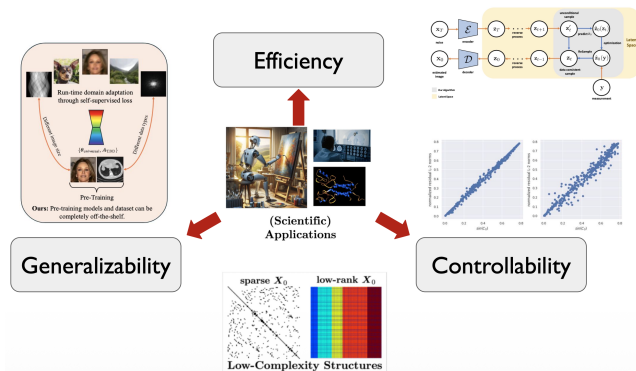
# S3: Understanding Low-Dimensional Structures in Diffusion Generative Models



Given enough data and ability to optimize: *inevitable emergence* of low-dim structures in trained deep networks!

Implications for efficiency, controllability, generalizability

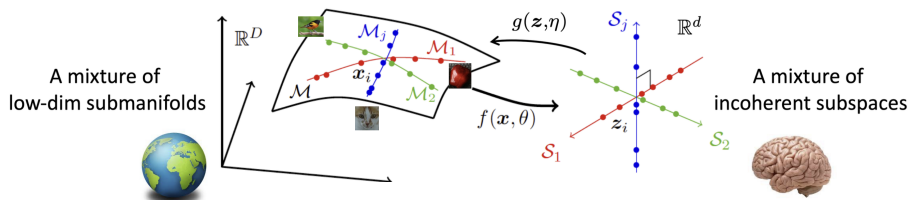
# S3: Understanding Low-Dimensional Structures in Diffusion Generative Models



Given enough data and ability to optimize: *inevitable emergence* of low-dim structures in trained deep networks!

Implications for efficiency, controllability, generalizability

# S4: Bottom-Up Understanding of Deep Networks for Vision



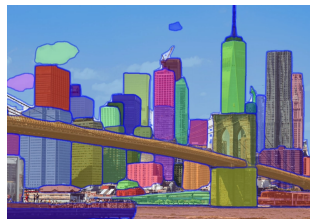
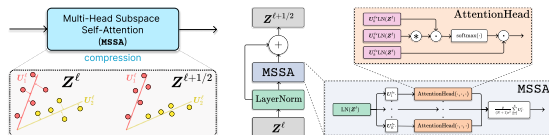
## The objective of learning:

*Identify* low-dim. distributions in sensed data of the world  
and *transform* to a **compact and structured** representation.

*Once we clarify this objective of learning, we can **design networks** to explicitly achieve these functions!*

# S4: Bottom-Up Understanding of Deep Networks for Vision

1. **Design**  $\varphi(z)$  s.t.  $z$  **optimal**  $\iff$  **good representation**
2. **Construct**  $f$  via **incremental optimization** of  $\varphi$
3. Learn any parameters of  $f$  from data



More *interpretable* (derivations!) and *less superfluous pieces*!

Thank You! Questions?

# Learning Deep Representations of Data Distributions

Sam Buchanan · Druv Pai · Peng Wang · Yi Ma

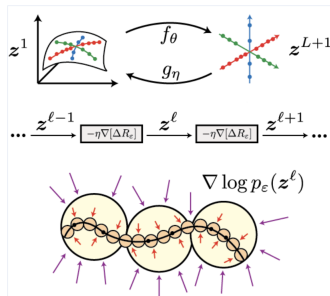
A modern fully open-source textbook exploring why and how deep neural networks learn compact and information-dense representations of high-dimensional real-world data.

```
@book{ldrdd2025,  
  title={Learning Deep Representations of Data Distributions},  
  author={Buchanan, Sam and Pai, Druv and Wang, Peng and Ma, Yi},  
  month=aug,  
  year={2025},  
  publisher={Online},  
  note={\url{https://ma-lab-berkeley.github.io/deep-representation-learning-book/.}}  
}
```

[Read the Book \(HTML\)](#)

[Read the Book \(PDF\)](#)

<https://ma-lab-berkeley.github.io/deep-representation-learning-book>



Version 1.0

Last Updated: October 17, 2025