

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

(ReduNet: Deep Networks from Maximizing Rate Reduction)

**Professor Yi Ma**

University of Hong Kong

October 19, 2025

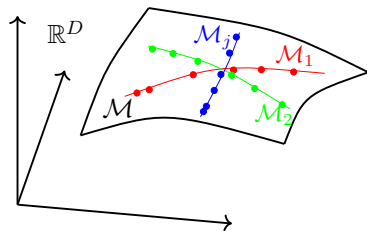
*“What I cannot create, I do not understand.”*  
– Richard Feynman

# Outline

- 1 Objectives for Learning from Data
  - Precursors and Motivations
  - Linear and Discriminative Representation (LDR)
- 2 Measure of Information Gain for Representations
  - Principle of Maximizing Coding Rate Reduction (MCR<sup>2</sup>)
  - Experimental Verification
- 3 White-Box Deep Networks from Optimizing Rate Reduction
  - Deep Networks as Projected Gradient Ascent
  - Convolution Networks from Shift Invariance
  - Experimental Results
  - Extension to White-Box Transformers via Sparse MCR<sup>2</sup>
- 4 Conclusions and Open Directions

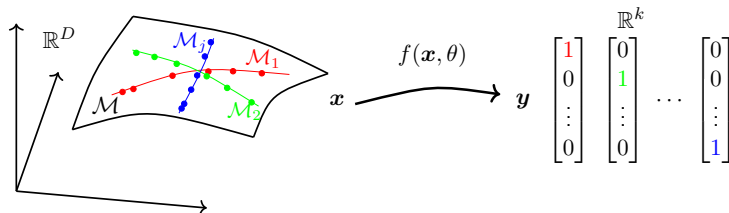
# High-Dim Data with Mixed **Nonlinear** Low-Dim Structures

**Figure: High-dimensional Real-World Data:** data samples  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m]$  in  $\mathbb{R}^D$  lying on a mixture of low-dimensional submanifolds  $\mathbf{X} \subset \cup_{j=1}^k \mathcal{M}_j \subset \mathbb{R}^D$ .



The main objective of learning from (samples of) such real-world data:  
**seek a most compact and structured representation of the data.**

# Fitting Class Labels via a Deep Network



**Figure: Black Box DNN for Classification:**  $y$  is the class label of  $x$  represented as a “one-hot” vector in  $\mathbb{R}^k$ . To learn a nonlinear mapping  $f(\cdot, \theta) : x \mapsto y$ , say modeled by a deep network, using cross-entropy (CE) loss.

$$\min_{\theta \in \Theta} \text{CE}(\theta, x, y) \doteq -\mathbb{E}[\langle y, \log[f(x, \theta)] \rangle] \approx -\frac{1}{m} \sum_{i=1}^m \langle y_i, \log[f(x_i, \theta)] \rangle. \quad (1)$$

*Prevalence of **neural collapse** during the terminal phase of deep learning training,*  
Papayan, Han, and Donoho, 2020.



# Fitting Class Labels via a Deep Network

In a supervised setting, using cross-entropy (CE) loss:

$$\min_{\theta \in \Theta} \text{CE}(\theta, \mathbf{x}, \mathbf{y}) \doteq -\mathbb{E}[\langle \mathbf{y}, \log[f(\mathbf{x}, \theta)] \rangle] \approx -\frac{1}{m} \sum_{i=1}^m \langle \mathbf{y}_i, \log[f(\mathbf{x}_i, \theta)] \rangle. \quad (2)$$

Issues (an elephant in the room):

- A large deep neural networks can **fit arbitrary data and labels**.
- Statistical and geometric meaning of internal features **not clear**.
- Task/data-dependent and **not robust nor truly invariant**.

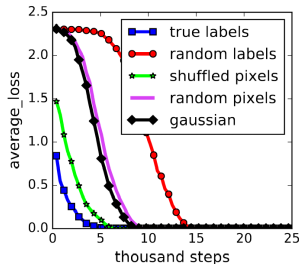


Figure: [Zhang et al, ICLR'17]

What did machines actually “learn” from doing this?

**In terms of interpolating, extrapolating, or representing the data?**

# A Hypothesis: Information Bottleneck

[Tishby & Zaslavsky, 2015]

A feature mapping  $f(\mathbf{x}, \theta)$  and a classifier  $g(\mathbf{z})$  trained for downstream classification:

$$\mathbf{x} \xrightarrow{f(\mathbf{x}, \theta)} \mathbf{z}(\theta) \xrightarrow{g(\mathbf{z})} \mathbf{y}.$$

**The IB Hypothesis:** Features learned in a deep network trying to

$$\max_{\theta \in \Theta} \text{IB}(\mathbf{x}, \mathbf{y}, \mathbf{z}(\theta)) \doteq I(\mathbf{z}(\theta), \mathbf{y}) - \beta I(\mathbf{x}, \mathbf{z}(\theta)), \quad \beta > 0, \quad (3)$$

where  $I(\mathbf{z}, \mathbf{y}) \doteq H(\mathbf{z}) - H(\mathbf{z}|\mathbf{y})$  and  $H(\mathbf{z})$  is the entropy of  $\mathbf{z}$ .

- **Minimal** informative features  $\mathbf{z}$  that most correlate with the label  $\mathbf{y}$
- Task and label-dependent, consequently sacrificing generalizability, robustness, or transferability

# Gap between Theory and Practice (a Bigger Elephant)

## For high-dimensional real data,

many statistical and information-theoretic concepts such as entropy, mutual information, K-L divergence, and maximum likelihood:

- curse of **dimensionality** for computation.
- ill-posed for **degenerate** distributions.
- lack guarantees with **finite** (or non-asymptotic) samples.

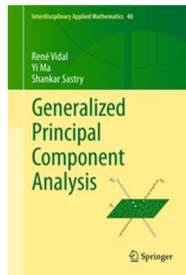
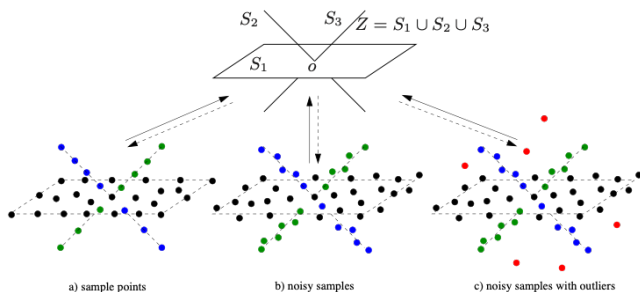
**Reality check:** principled formulations are replaced with approximate bounds, grossly simplifying assumptions, heuristics, even *ad hoc* tricks and hacks.

**How to provide any performance guarantees at all?**

# A Principled Computational Approach

For high-dim data with mixed **low-dim linear/Gaussian** structures:

**learn to compress, and compress to learn!**



**Generalized PCA** for mixture of subspaces [Vidal, Ma, and Sastry, 2005]

# Clustering (or Classification) via Compression

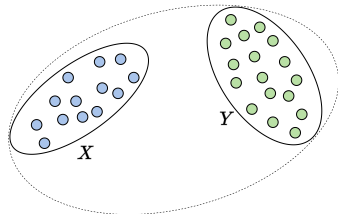
[Yi Ma, Harm Derksen, Wei Hong, and John Wright, TPAMI'07]

## A Fundamental Idea:

Data belong to mixed low-dim structures should be compressible.

## Cluster Criterion:

Whether the number of binary bits required to store the data is less (information gain):



$$\#bits(\mathbf{X} \cup \mathbf{Y}) \geq \#bits(\mathbf{X}) + \#bits(\mathbf{Y})?$$

*"The whole is greater than the sum of the parts."*  
– Aristotle, 320 BC

# Coding Length Function for Subspace-Like Data

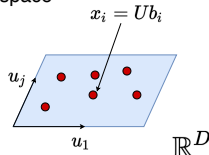
## Theorem (Ma, TPAMI'07)

The number of bits needed to encode data  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m] \in \mathbb{R}^{D \times m}$  up to a precision  $\|\mathbf{x} - \hat{\mathbf{x}}\|_2 \leq \epsilon$  is bounded by:

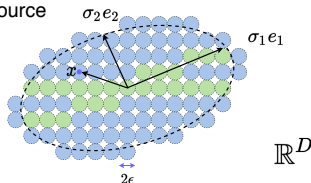
$$L(\mathbf{X}, \epsilon) \doteq \left( \frac{m + D}{2} \right) \log \det \left( \mathbf{I} + \frac{D}{m\epsilon^2} \mathbf{X} \mathbf{X}^\top \right).$$

This can be derived from constructively quantifying SVD of  $\mathbf{X}$  or by sphere packing  $\text{vol}(\mathbf{X})$  as samples of a noisy Gaussian source.

Linear subspace



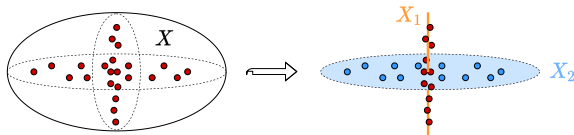
Gaussian source



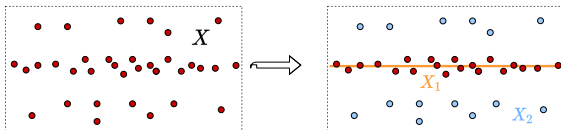
# Cluster to Compress

$$L(\mathbf{X}) \geq L^c(\mathbf{X}) \doteq L(\mathbf{X}_1) + L(\mathbf{X}_2) + H(|\mathbf{X}_1|, |\mathbf{X}_2|)?$$

partitioning:



sifting:



# A Greedy Algorithm

Seek a partition of the data  $\mathbf{X} \rightarrow [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_k]$  such that

$$\min L^c(\mathbf{X}) \doteq L(\mathbf{X}_1) + \dots + L(\mathbf{X}_k) + H(|\mathbf{X}_1|, \dots, |\mathbf{X}_k|).$$

Optimize with a *bottom-up pair-wise* merging algorithm [Ma, TPAMI'07]:

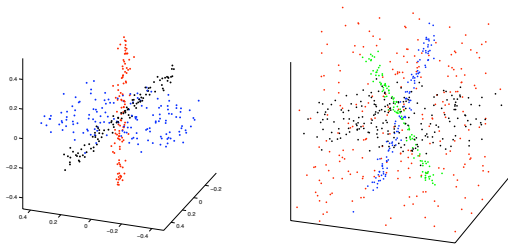
- 1: **input:** the data  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m] \in \mathbb{R}^{D \times m}$  and a distortion  $\epsilon^2 > 0$ .
- 2: initialize  $\mathcal{S}$  as a set of sets with a single datum  $\{S = \{\mathbf{x}\} \mid \mathbf{x} \in \mathbf{X}\}$ .
- 3: **while**  $|\mathcal{S}| > 1$  **do**
- 4: choose distinct sets  $S_1, S_2 \in \mathcal{S}$  such that  

$$L^c(S_1 \cup S_2) - L^c(S_1, S_2)$$
 is minimal.
- 5: **if**  $L^c(S_1 \cup S_2) - L^c(S_1, S_2) \geq 0$  **then** break;
- 6: **else**  $\mathcal{S} := (\mathcal{S} \setminus \{S_1, S_2\}) \cup \{S_1 \cup S_2\}$ .
- 7: **end**
- 8: **output:**  $\mathcal{S}$

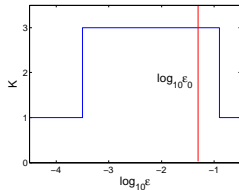
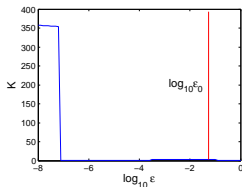


# Surprisingly Good Performance

Empirically, **find global optimum** and **extremely robust to outliers**

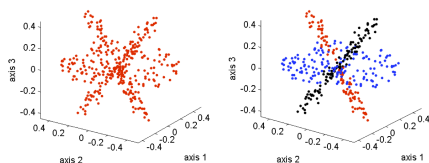


A strikingly sharp **phase transition** w.r.t. quantization  $\epsilon$

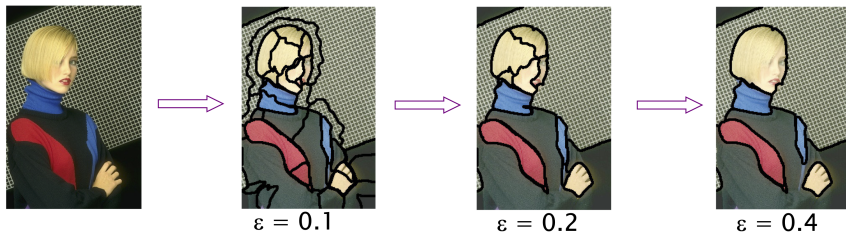


# Clustering by Minimizing Coding Length

*Segmentation of Multivariate Mixed Data via Lossy Coding and Compression,*  
Yi Ma et. al., TPAMI, 2007.



**State of the art unsupervised image segmentation (IJCV 2011):**



# Natural Image Segmentation [Mobahi et.al., IJCV'09]

**Compression alone**, without any supervision, leads to **state of the art** segmentation on natural images (and many other types of data).



(a) Animals



(b) Buildings



(c) Landscape



(d) People



(e) Water

# Represent Multi-class Multi-dimensional Data

Given samples

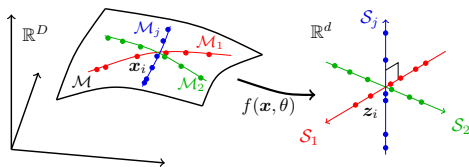
$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m] \subset \cup_{j=1}^k \mathcal{M}_j,$$

**seek a good representation**

$$\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_m] \subset \mathbb{R}^d$$

through a continuous mapping:

$$f(\mathbf{x}, \theta) : \mathbf{x} \in \mathbb{R}^D \mapsto \mathbf{z} \in \mathbb{R}^d.$$



Goals of “**re-present**” the data:

- **compression**: from high-dimensional samples to compact features.
- **linearization**: from nonlinear structures  $\cup_{j=1}^k \mathcal{M}_j$  to linear  $\cup_{j=1}^k \mathcal{S}_j$ .
- **sparsity**: from separable components  $\mathcal{M}_j$ 's to incoherent  $\mathcal{S}_j$ 's.

# Seeking a Linear Discriminative Representation (LDR)

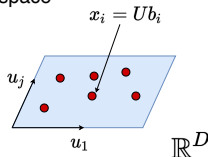
**Desiderata:** Representation  $z = f(x, \theta)$  have the following properties:

- ① *Within-Class Compressible:* Features of the same class/cluster should be highly compressed in a **low-dimensional** linear subspace.
- ② *Between-Class Discriminative:* Features of different classes/clusters should be in highly **incoherent** linear subspaces.
- ③ *Maximally Informative:* Dimension (or variance) of features for each class/cluster should be **the same as that of the data**.

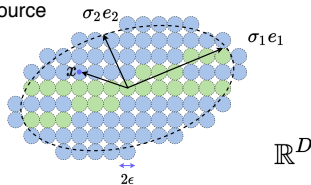
**Is there a principled objective for all such properties, together?**

# Compactness Measure for Linear/Gaussian Representation

Linear subspace



Gaussian source



## Theorem (Coding Length, Ma & Derksen TPAMI'07)

The number of bits needed to encode data  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m] \in \mathbb{R}^{D \times m}$  up to a precision  $\|\mathbf{x} - \hat{\mathbf{x}}\|_2 \leq \epsilon$  is bounded by:

$$L(\mathbf{X}, \epsilon) \doteq \left( \frac{m + D}{2} \right) \log \det \left( \mathbf{I} + \frac{D}{m\epsilon^2} \mathbf{X} \mathbf{X}^\top \right).$$

This can be derived from constructively quantifying SVD of  $\mathbf{X}$  or by sphere packing  $\text{vol}(\mathbf{X})$  as samples of a noisy Gaussian source.

# Compactness Measure for Linear/Gaussian Representation

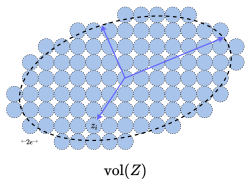
If  $\mathbf{X}$  is not (piecewise) linear or Gaussian, consider a **nonlinear** mapping:

$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m] \in \mathbb{R}^{D \times m} \xrightarrow{f(\mathbf{x}, \theta)} \mathbf{Z}(\theta) = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_m] \in \mathbb{R}^{d \times m}.$$

The average coding length per sample (rate) subject to a distortion  $\epsilon$ :

$$R(\mathbf{Z}, \epsilon) \doteq \frac{1}{2} \log \det \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z} \mathbf{Z}^\top \right). \quad (4)$$

**Rate distortion is an intrinsic measure for the volume of all features.**



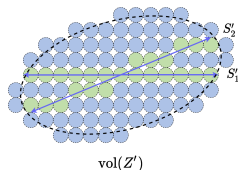
# Compactness Measure for Mixed Linear Representations

The features  $\mathbf{Z}$  of **multi-class** data

$$\mathbf{X} = \mathbf{X}_1 \cup \mathbf{X}_2 \cup \cdots \cup \mathbf{X}_k \subset \cup_{j=1}^k \mathcal{M}_j.$$

may be partitioned into **multiple** subsets:

$$\mathbf{Z} = \mathbf{Z}_1 \cup \mathbf{Z}_2 \cup \cdots \cup \mathbf{Z}_k \subset \cup_{j=1}^k \mathcal{S}_j.$$



W.r.t. this partition, the **average coding rate** is:

$$R^c(\mathbf{Z}, \epsilon \mid \mathbf{\Pi}) \doteq \sum_{j=1}^k \frac{\text{tr}(\mathbf{\Pi}_j)}{2m} \log \det \left( \mathbf{I} + \frac{d}{\text{tr}(\mathbf{\Pi}_j) \epsilon^2} \mathbf{Z} \mathbf{\Pi}_j \mathbf{Z}^\top \right), \quad (5)$$

where  $\mathbf{\Pi} = \{\mathbf{\Pi}_j \in \mathbb{R}^{m \times m}\}_{j=1}^k$  encode the membership of the  $m$  samples in the  $k$  classes: the diagonal entry  $\mathbf{\Pi}_j(i, i)$  of  $\mathbf{\Pi}_j$  is the probability of sample  $i$  belonging to subset  $j$ .  $\Omega \doteq \{\mathbf{\Pi} \mid \sum \mathbf{\Pi}_j = \mathbf{I}, \mathbf{\Pi}_j \geq \mathbf{0}\}$



# Measure for Linear Discriminative Representation (LDR)

**A fundamental idea:** maximize the **difference** between the coding rate of all features and the average rate of features in each of the classes:

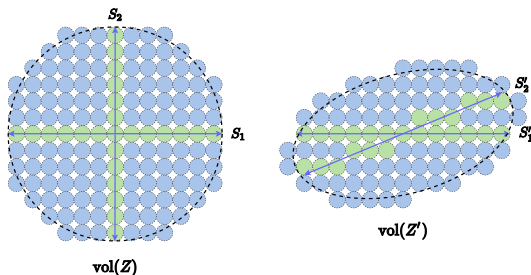
$$\Delta R(\mathbf{Z}, \mathbf{\Pi}, \epsilon) = \underbrace{\frac{1}{2} \log \det \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z} \mathbf{Z}^\top \right)}_R - \underbrace{\sum_{j=1}^k \frac{\text{tr}(\mathbf{\Pi}_j)}{2m} \log \det \left( \mathbf{I} + \frac{d}{\text{tr}(\mathbf{\Pi}_j)\epsilon^2} \mathbf{Z} \mathbf{\Pi}_j \mathbf{Z}^\top \right)}_{R^c}.$$

This difference is called **rate reduction** (measuring **information gain**):

- Large  $R$ : **expand** all features  $\mathbf{Z}$  as **large** as possible.
- Small  $R^c$ : **compress** each class  $\mathbf{Z}_j$  as **small** as possible.

**Slogan: similarity contracts and dissimilarity contrasts!**

# Interpretation of MCR<sup>2</sup>: Sphere Packing and Counting



**Example:** two subspaces  $S_1$  and  $S_2$  in  $\mathbb{R}^2$ .

- $\log \#(\text{green spheres} + \text{blue spheres}) = \text{rate of span of all samples } R.$
- $\log \#(\text{green spheres}) = \text{rate of the two subspaces } R^c.$
- $\log \#(\text{blue spheres}) = \text{rate reduction } \Delta R.$

# Comparison to Contrastive Learning

[Hadsell, Chopra, and LeCun, CVPR'06]

When  $k$  is large, a randomly chosen **pair**  $(x_i, x_j)$  is of high probability belonging to different classes. Minimize the **contrastive loss**:

$$\min -\log \frac{\exp(\langle z_i, z'_i \rangle)}{\sum_{j \neq i} \exp(\langle z_i, z_j \rangle)}.$$

The learned features of such pairs of samples together with their augmentations  $\mathbf{Z}_i$  and  $\mathbf{Z}_j$  should have large rate reduction:

$$\max_{ij} \sum \Delta R_{ij} \doteq R(\mathbf{Z}_i \cup \mathbf{Z}_j, \epsilon) - \frac{1}{2}(R(\mathbf{Z}_i, \epsilon) + R(\mathbf{Z}_j, \epsilon)).$$

**MCR<sup>2</sup> contrasts triplets, quadruplets, or any number of sets.**

# Principle of Maximal Coding Rate Reduction (MCR<sup>2</sup>)

[Yu, Chan, You, Song, Ma, NeurIPS2020]

Learn a mapping  $f(\mathbf{x}, \theta)$  (for a given partition  $\Pi$ ):

$$\mathbf{X} \xrightarrow{f(\mathbf{x}, \theta)} \mathbf{Z}(\theta) \xrightarrow{\Pi, \epsilon} \Delta R(\mathbf{Z}(\theta), \Pi, \epsilon) \quad (6)$$

so as to **Maximize the Coding Rate Reduction (MCR<sup>2</sup>)**:

$$\begin{aligned} \max_{\theta} \quad & \Delta R(\mathbf{Z}(\theta), \Pi, \epsilon) = R(\mathbf{Z}(\theta), \epsilon) - R^c(\mathbf{Z}(\theta), \epsilon \mid \Pi), \\ \text{subject to} \quad & \|\mathbf{Z}_j(\theta)\|_F^2 = m_j, \Pi \in \Omega. \end{aligned} \quad (7)$$

Since  $\Delta R$  is *monotonic* in the scale of  $\mathbf{Z}$ , one needs to:

**normalize the features  $\mathbf{z} = f(\mathbf{x}, \theta)$  so as to compare  $\mathbf{Z}(\theta)$  and  $\mathbf{Z}(\theta')$ !**

Batch normalization, Sergey Ioffe and Christian Szegedy, 2015.

Layer normalization'16, instance normalization'16; group normalization'18...

# Theoretical Justification of the MCR<sup>2</sup> Principle

## Theorem (Informal Statement [Yu et.al., NeurIPS2020])

*Suppose  $\mathbf{Z}^* = \mathbf{Z}_1^* \cup \dots \cup \mathbf{Z}_k^*$  is the optimal solution that maximizes the rate reduction (7). We have:*

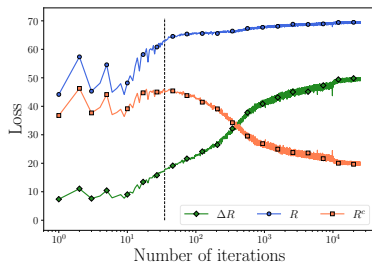
- *Between-class Discriminative: As long as the ambient space is adequately large ( $d \geq \sum_{j=1}^k d_j$ ), the subspaces are all orthogonal to each other, i.e.  $(\mathbf{Z}_i^*)^\top \mathbf{Z}_j^* = \mathbf{0}$  for  $i \neq j$ .*
- *Maximally Informative Representation: As long as the coding precision is adequately high, i.e.,  $\epsilon^4 < \min_j \left\{ \frac{m_j}{m} \frac{d^2}{d_j^2} \right\}$ , each subspace achieves its maximal dimension, i.e.  $\text{rank}(\mathbf{Z}_j^*) = d_j$ . In addition, the largest  $d_j - 1$  singular values of  $\mathbf{Z}_j^*$  are equal.*

A new slogan, beyond Aristotle:

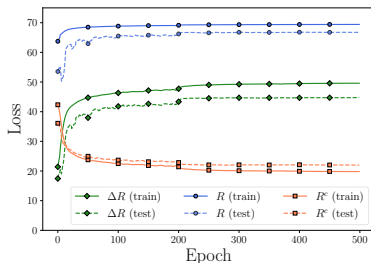
**The whole is to be maximally greater than the sum of the parts!**

# Experiment I: Supervised Deep Learning

**Experimental Setup:** Train  $f(x, \theta)$  as ResNet18 on the CIFAR10 dataset, feature  $z$  dimension  $d = 128$ , precision  $\epsilon^2 = 0.5$ .



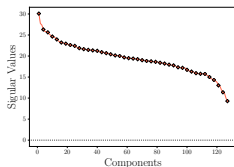
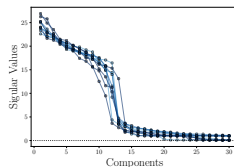
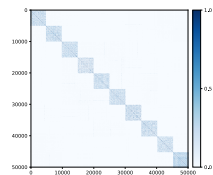
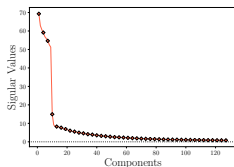
(a)



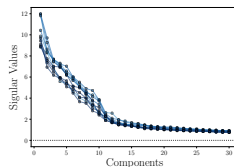
(b)

**Figure:** (a). Evolution of  $R, R^c, \Delta R$  during the training process; (b). Training loss versus testing loss.

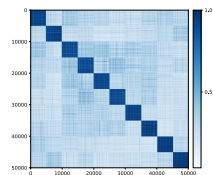
# Visualization of Learned Representations $\mathbb{Z}$

(a)  $\text{MCR}^2$  (overall)(b)  $\text{MCR}^2$  (PCA of every class)(c)  $\text{MCR}^2$  (cosine similarity)

(d) CE (overall)



(e) CE (PCA of every class)

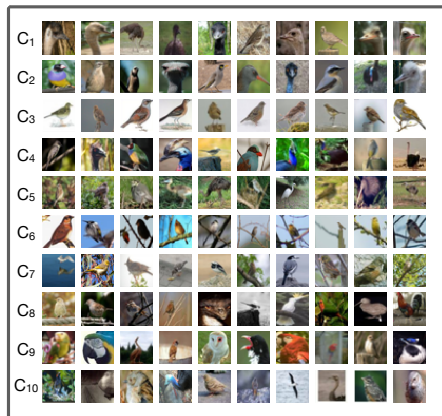


(f) CE (cosine similarity)

Figure: PCA of learned representations from  $\text{MCR}^2$  and cross-entropy.

**No neural collapse!**

# Visualization - Samples along Principal Components



(a) Bird



(b) Ship

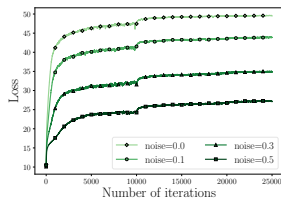
**Figure:** Top-10 “principal” images for class - “Bird” and “Ship” in the CIFAR10.



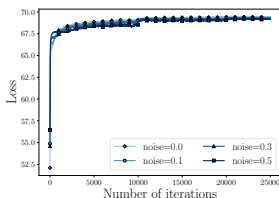
# Experiment II: Robustness to Label Noise

	RATIO=0.0	RATIO=0.1	RATIO=0.2	RATIO=0.3	RATIO=0.4	RATIO=0.5
CE TRAINING	0.939	0.909	0.861	0.791	0.724	0.603
MCR <sup>2</sup> TRAINING	<b>0.940</b>	<b>0.911</b>	<b>0.897</b>	<b>0.881</b>	<b>0.866</b>	<b>0.843</b>

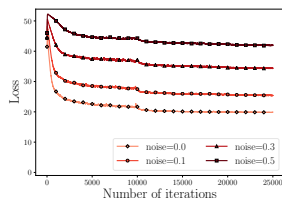
Table 1: Classification results with features learned with labels corrupted at different levels.



(a)  $\Delta R(\mathbf{Z}(\theta), \Pi, \epsilon)$



(b)  $R(\mathbf{Z}(\theta), \epsilon)$



(c)  $R^c(\mathbf{Z}(\theta), \epsilon | \Pi)$

**Figure:** Evolution of  $R, R^c, \Delta R$  of MCR<sup>2</sup> during training with corrupted labels.

**Represent only what can be jointly compressed.**

# Deep Networks from Optimizing Rate Reduction

$$\mathbf{X} \xrightarrow{f(\mathbf{x}, \theta)} \mathbf{Z}(\theta); \quad \max_{\theta} \Delta R(\mathbf{Z}(\theta), \mathbf{\Pi}, \epsilon).$$

Final features learned by MCR<sup>2</sup> are more interpretable and robust, **but**:

- The borrowed deep network (e.g. ResNet) is still a “black box”!
- Why is a “deep” architecture necessary, and how wide and deep?
- What are the roles of the “linear and nonlinear” operators?
- Why “multi-channel” convolutions?
- ...

**Replace black box networks with entirely “white box” networks?**

# Projected Gradient Ascent for Rate Reduction

Recall the rate reduction objective:

$$\max_{\mathbf{Z}} \Delta R(\mathbf{Z}) \doteq \underbrace{\frac{1}{2} \log \det \left( \mathbf{I} + \alpha \mathbf{Z} \mathbf{Z}^* \right)}_{R(\mathbf{Z})} - \underbrace{\sum_{j=1}^k \frac{\gamma_j}{2} \log \det \left( \mathbf{I} + \alpha_j \mathbf{Z} \mathbf{\Pi}^j \mathbf{Z}^* \right)}_{R_c(\mathbf{Z}, \mathbf{\Pi})}, \quad (8)$$

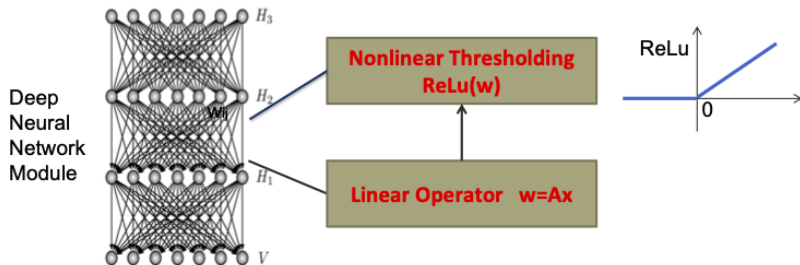
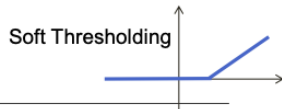
where  $\alpha = d/(m\epsilon^2)$ ,  $\alpha_j = d/(\text{tr}(\mathbf{\Pi}^j)\epsilon^2)$ ,  $\gamma_j = \text{tr}(\mathbf{\Pi}^j)/m$  for  $j = 1, \dots, k$ .

Consider directly maximizing  $\Delta R$  with **projected gradient ascent** (PGA):

$$\mathbf{Z}_{\ell+1} \propto \mathbf{Z}_{\ell} + \eta \cdot \left. \frac{\partial \Delta R}{\partial \mathbf{Z}} \right|_{\mathbf{Z}_{\ell}} \quad \text{subject to} \quad \mathbf{Z}_{\ell+1} \subset \mathbb{S}^{d-1}. \quad (9)$$

ISTA: Sparse Recovery via  $\ell^1$  (Wright and Ma, 2022)**CONTEXT – Basic algorithm (ISTA)****Algorithm 8.1** Iterative Soft-Thresholding Algorithm (ISTA) for BPDN

- 1: **Problem:**  $\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$ , given  $\mathbf{y} \in \mathbb{R}^d$ ,  $\mathbf{A} \in \mathbb{R}^{d \times n}$ .
- 2: **Input:**  $\mathbf{x}_0 \in \mathbb{R}^n$  and  $L \geq \lambda_{\max}(\mathbf{A}^T \mathbf{A})$ .
- 3: **while**  $\mathbf{x}_k$  not converged ( $k = 1, 2, \dots$ ) **do**
- 4:    $\mathbf{w}_k \leftarrow \mathbf{x}_k - \frac{1}{L} \mathbf{A}^T (\mathbf{A}\mathbf{x}_k - \mathbf{y})$ .
- 5:    $\mathbf{x}_{k+1} \leftarrow \text{soft}(\mathbf{w}_k, \lambda/L)$ .
- 6: **end while**
- 7: **Output:**  $\mathbf{x}_* \leftarrow \mathbf{x}_k$ .



# Learned ISTA (Gregor and LeCun, ICML 2010)

## CONTEXT – Learned ISTA (LISTA)

If only interested in one instance:  $y = Ax$  AND with many training data:  $\{(y_i, x_i)\}$ .  
We can **optimize the optimization path** of ISTA using supervised learning:

---

### Algorithm 3 LISTA::fprop

---

```

LISTA :: fprop( $X, Z, W_e, S, \theta$ )
;; Arguments are passed by reference.
;; variables  $Z(t)$ ,  $C(t)$  and  $B$  are saved for bprop.
 $B = W_e X$ ;  $Z(0) = h_\theta(B)$ 
for  $t = 1$  to  $T$  do
     $C(t) = B + SZ(t-1)$ 
     $Z(t) = h_\theta(C(t))$ 
end for
 $Z = Z(T)$ 

```

---



---

### Algorithm 4 LISTA::bprop

---

```

LISTA :: bprop( $Z^*, X, Z, W_e, S, \theta, \delta X, \delta W_e, \delta S, \delta \theta$ )
;; Arguments are passed by reference.
;; Variables  $Z(t)$ ,  $C(t)$ , and  $B$  were saved in fprop.
Initialize:  $\delta B = 0$ ;  $\delta S = 0$ ;  $\delta \theta = 0$ 
 $\delta Z(T) = (Z(T) - Z^*)$ 
for  $t = T$  down to  $1$  do
     $\delta C(t) = h'_\theta(C(t)).\delta Z(t)$ 
     $\delta \theta = \delta \theta - \text{sign}(C(t)).\delta C(t)$ 
     $\delta B = \delta B + \delta C(t)$ 
     $\delta S = \delta S + \delta C(t)Z(t-1)^T$ 
     $\delta Z(t-1) = S^T \delta C(t)$ 
end for
 $\delta B = \delta B + h'_\theta(B).\delta Z(0)$ 
 $\delta \theta = \delta \theta - \text{sign}(B).h'_\theta(B)\delta Z(0)$ 
 $\delta W_e = \delta B X^T$ ;  $\delta X = W_e^T \delta B$ 

```

---

Learning fast approximations of sparse coding, K. Gregor and Y. LeCun, ICML 2010.

# Gradients of the Two Terms

The derivatives  $\frac{\partial R(\mathbf{Z})}{\partial \mathbf{Z}}$  and  $\frac{\partial R_c(\mathbf{Z}, \mathbf{\Pi})}{\partial \mathbf{Z}}$  are:

$$\left. \frac{1}{2} \frac{\partial \log \det(\mathbf{I} + \alpha \mathbf{Z} \mathbf{Z}^*)}{\partial \mathbf{Z}} \right|_{\mathbf{Z}_\ell} = \underbrace{\alpha (\mathbf{I} + \alpha \mathbf{Z}_\ell \mathbf{Z}_\ell^*)^{-1}}_{\mathbf{E}_\ell \in \mathbb{R}^{d \times d}} \mathbf{Z}_\ell, \quad (10)$$

$$\left. \frac{1}{2} \frac{\partial (\gamma_j \log \det(\mathbf{I} + \alpha_j \mathbf{Z} \mathbf{\Pi}^j \mathbf{Z}^*))}{\partial \mathbf{Z}} \right|_{\mathbf{Z}_\ell} = \gamma_j \underbrace{\alpha_j (\mathbf{I} + \alpha_j \mathbf{Z}_\ell \mathbf{\Pi}^j \mathbf{Z}_\ell^*)^{-1}}_{\mathbf{C}_\ell^j \in \mathbb{R}^{d \times d}} \mathbf{Z}_\ell \mathbf{\Pi}^j. \quad (11)$$

Hence the gradient  $\frac{\partial \Delta R(\mathbf{Z})}{\partial \mathbf{Z}}$  is:

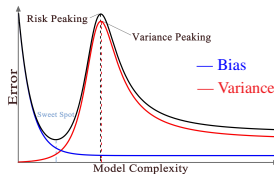
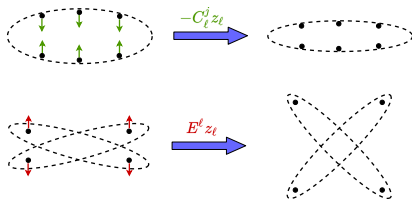
$$\left. \frac{\partial \Delta R}{\partial \mathbf{Z}} \right|_{\mathbf{Z}_\ell} = \underbrace{\mathbf{E}_\ell}_{\text{Expansion}} \mathbf{Z}_\ell - \sum_{j=1}^k \gamma_j \underbrace{\mathbf{C}_\ell^j}_{\text{Compression}} \mathbf{Z}_\ell \mathbf{\Pi}^j \in \mathbb{R}^{d \times m}. \quad (12)$$

# Interpretation of the Linear Operators $E$ and $C^j$

For any  $z_\ell \in \mathbb{R}^d$ , we have

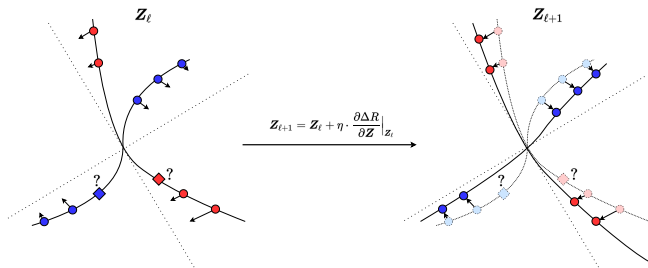
$$E_\ell z_\ell = \alpha(z_\ell - Z_\ell q_\ell^*) \quad \text{with} \quad q_\ell^* \doteq \arg \min_{q_\ell} \alpha \|z_\ell - Z_\ell q_\ell\|_2^2 + \|q_\ell\|_2^2.$$

$E_\ell z_\ell$  and  $C_\ell^j z_\ell$  are the “residuals” of  $z_\ell$  against the subspaces spanned by columns of  $Z_\ell$  and  $Z_\ell^j$ , respectively.



Such “auto” ridge regressions **do not overfit** even with redundant random regressors, due to a “double descent” risk [Yang, ICML'20]!

# Incremental Deformation via Gradient Flow



Extrapolate the gradient  $\frac{\partial \Delta R(\mathbf{Z})}{\partial \mathbf{Z}}$  from training samples  $\mathbf{Z}$  to all  $\mathbf{z} \in \mathbb{R}^d$ :

$$\frac{\partial \Delta R}{\partial \mathbf{Z}} \Big|_{\mathbf{Z}_\ell} = \mathbf{E}_\ell \mathbf{Z}_\ell - \sum_{j=1}^k \gamma_j \mathbf{C}_\ell^j \mathbf{Z}_\ell \underbrace{\mathbf{\Pi}^j}_{\text{known}} \in \mathbb{R}^{d \times m}, \quad (13)$$

$$g(\mathbf{z}_\ell, \boldsymbol{\theta}_\ell) \doteq \mathbf{E}_\ell \mathbf{z}_\ell - \sum_{j=1}^k \gamma_j \mathbf{C}_\ell^j \mathbf{z}_\ell \underbrace{\boldsymbol{\pi}^j(\mathbf{z}_\ell)}_{\text{unknown}} \in \mathbb{R}^d. \quad (14)$$



## Estimate of the Membership $\pi^j(z_\ell)$

Estimate the membership  $\pi^j(z_\ell)$  with “softmax” on the residuals  $\|C_\ell^j z_\ell\|$ :

$$\pi^j(z_\ell) \approx \hat{\pi}^j(z_\ell) \doteq \frac{\exp(-\lambda \|C_\ell^j z_\ell\|)}{\sum_{j=1}^k \exp(-\lambda \|C_\ell^j z_\ell\|)} \in [0, 1]. \quad (15)$$

Thus the weighted residuals for contracting:

$$\sigma\left([C_\ell^1 z_\ell, \dots, C_\ell^k z_\ell]\right) \doteq \sum_{j=1}^k \gamma_j C_\ell^j z_\ell \cdot \hat{\pi}^j(z_\ell) \in \mathbb{R}^d. \quad (16)$$

Many alternatives, e.g. enforcing all features to be in the first quadrant:

$$\sigma(z_\ell) \approx z_\ell - \sum_{j=1}^k \text{ReLU}(P_\ell^j z_\ell), \quad (17)$$

# The ReduNet for Optimizing Rate Reduction

Iterative projected gradient ascent (PGA) :

$$\mathbf{z}_{\ell+1} \propto \mathbf{z}_{\ell} + \eta \cdot \underbrace{\left[ \mathbf{E}_{\ell} \mathbf{z}_{\ell} + \sigma([C_{\ell}^1 \mathbf{z}_{\ell}, \dots, C_{\ell}^k \mathbf{z}_{\ell}]) \right]}_{g(\mathbf{z}_{\ell}, \boldsymbol{\theta}_{\ell})} \quad \text{s.t.} \quad \mathbf{z}_{\ell+1} \in \mathbb{S}^{d-1}, \quad (18)$$

$f(\mathbf{x}, \boldsymbol{\theta}) = \phi^L \circ \phi^{L-1} \circ \dots \circ \phi^0(\mathbf{x})$ , with  $\phi^{\ell}(\mathbf{z}_{\ell}, \boldsymbol{\theta}_{\ell}) \doteq \mathcal{P}_{\mathbb{S}^{d-1}}[\mathbf{z}_{\ell} + \eta \cdot g(\mathbf{z}_{\ell}, \boldsymbol{\theta}_{\ell})]$ .

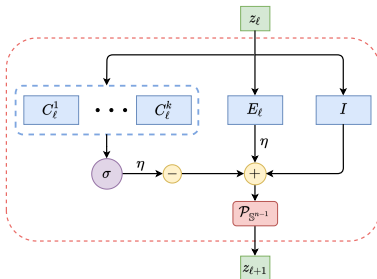


Figure: One layer of the **ReduNet**: one PGA iteration.

# The ReduNet versus ResNet or ResNeXt

Iterative projected gradient ascent (PGA):

$$\mathbf{z}_{\ell+1} \propto \mathbf{z}_{\ell} + \eta \cdot \underbrace{\left[ \mathbf{E}_{\ell} \mathbf{z}_{\ell} + \sigma([C_{\ell}^1 \mathbf{z}_{\ell}, \dots, C_{\ell}^k \mathbf{z}_{\ell}]) \right]}_{g(\mathbf{z}_{\ell}, \theta_{\ell})} \quad \text{s.t.} \quad \mathbf{z}_{\ell+1} \in \mathbb{S}^{d-1}, \quad (19)$$

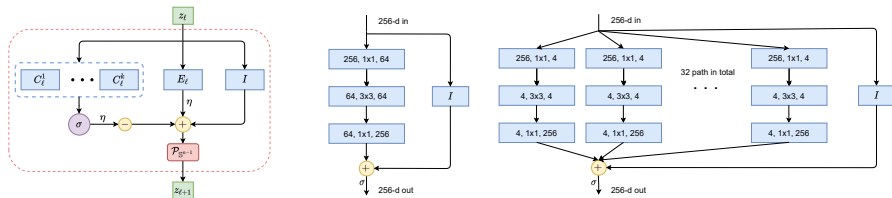


Figure: Left: **ReduNet**. Middle and Right: **ResNet** [He et. al. 2016] and **ResNeXt** [Xie et. al. 2017] (hundreds of layers).

**Forward construction instead of back propagation!**<sup>1</sup>

<sup>1</sup> *The Forward-Forward Algorithm*, G. Hinton, 2022.

# The ReduNet versus Mixture of Experts

Approximate iterative projected gradient ascent (PGA) :

$$z_{\ell+1} \propto z_{\ell} + \eta \cdot \underbrace{\left[ E_{\ell} z_{\ell} + \sigma([C_{\ell}^1 z_{\ell}, \dots, C_{\ell}^k z_{\ell}]) \right]}_{g(z_{\ell}, \theta_{\ell})} \quad \text{s.t.} \quad z_{\ell+1} \in \mathbb{S}^{d-1}, \quad (20)$$

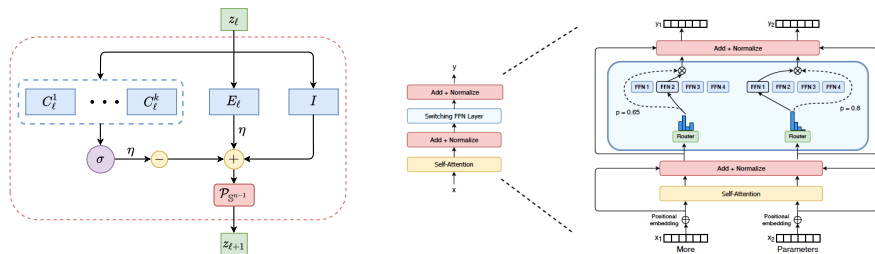


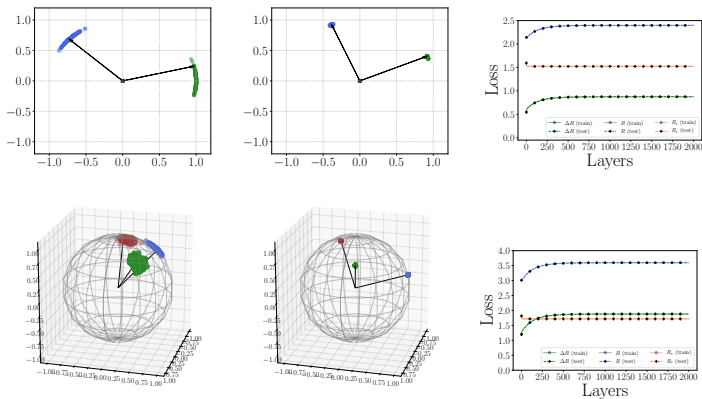
Figure: Left: ReduNet layer. Right: Mixture of Experts [Shazeer et. al. 2017] or Switched Transformer [Fedus et. al. 2021] (1.7 trillion parameters).

**Forward construction instead of back propagation!<sup>2</sup>**

<sup>2</sup> The Forward-Forward Algorithm, G. Hinton, 2022.

# ReduNet Features for Mixture of Gaussians

$L = 2000$ -Layers ReduNet:  $m = 500, \eta = 0.5, \epsilon = 0.1$ .



**Figure:** Left: original samples  $X$  and ReduNet features  $Z = f(Z, \theta)$  for 2D and 3D Mixture of Gaussians. Right: plots for the progression of values of the rates.

# Group Invariant Classification

Feature mapping  $f(x, \theta)$  is invariant to a group of transformations:

$$\text{Group Invariance: } f(x \circ g, \theta) \sim f(x, \theta), \quad \forall g \in \mathbb{G}, \quad (21)$$

where “ $\sim$ ” indicates two features belonging to the same equivalent class.

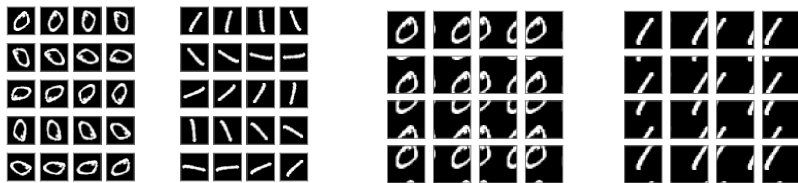


Figure: Left: 1D rotation  $\mathbb{S}^1$ ; Right: 2D cyclic translation  $\mathcal{T}^2$ .

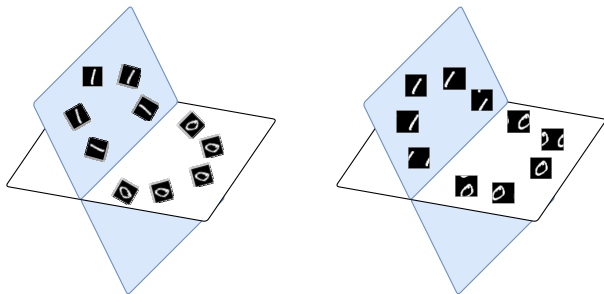
1. Fooling CNNs with simple transformations, Engstrom et.al., 2017.
2. Why do deep convolutional networks generalize so poorly to small image transformations? Azulay & Weiss, 2018.

# Group Invariant Classification

Feature mapping  $f(x, \theta)$  is invariant to a group of transformations:

$$\text{Group Invariance: } f(x \circ g, \theta) \sim f(x, \theta), \quad \forall g \in \mathbb{G}, \quad (22)$$

where “ $\sim$ ” indicates two features belonging to the same equivalent class.



**Figure:** Embed all equivariant samples to the same subspace.

# Circulant Matrix and Convolution

Given a vector  $\mathbf{z} = [z_0, z_1, \dots, z_{n-1}]^* \in \mathbb{R}^n$ , we may arrange all its circular shifted versions in a circulant matrix form as

$$\text{circ}(\mathbf{z}) \doteq \begin{bmatrix} z_0 & z_{n-1} & \dots & z_2 & z_1 \\ z_1 & z_0 & z_{n-1} & \dots & z_2 \\ \vdots & z_1 & z_0 & \ddots & \vdots \\ z_{n-2} & \vdots & \ddots & \ddots & z_{n-1} \\ z_{n-1} & z_{n-2} & \dots & z_1 & z_0 \end{bmatrix} \in \mathbb{R}^{n \times n}. \quad (23)$$

A circular (or cyclic) convolution:

$$\text{circ}(\mathbf{z}) \cdot \mathbf{x} = \mathbf{z} \circledast \mathbf{x}, \quad \text{where} \quad (\mathbf{z} \circledast \mathbf{x})_i = \sum_{j=0}^{n-1} x_j z_{i+n-j \bmod n}. \quad (24)$$



# Convolutions from Cyclic Shift Invariance

Given a set of sample vectors  $\mathbf{Z} = [\mathbf{z}^1, \dots, \mathbf{z}^m]$ , construct the ReduNet from cyclic-shift augmented families  $\mathbf{Z} = [\text{circ}(\mathbf{z}^1), \dots, \text{circ}(\mathbf{z}^m)]$ .

## Proposition (Convolution Structures of $\mathbf{E}$ and $\mathbf{C}^j$ )

*The linear operator in the ReduNet:*

$$\mathbf{E} = \alpha \left( \mathbf{I} + \alpha \sum_{i=1}^m \text{circ}(\mathbf{z}^i) \text{circ}(\mathbf{z}^i)^* \right)^{-1}$$

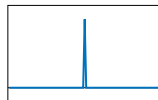
*is a circulant matrix and represents a circular convolution:*

$$\mathbf{E}\mathbf{z} = \mathbf{e} \circledast \mathbf{z},$$

*where  $\mathbf{e}$  is the first column vector of  $\mathbf{E}$ . Similarly, the operators  $\mathbf{C}^j$  associated with subsets  $\mathbf{Z}^j$  are also circular convolutions.*

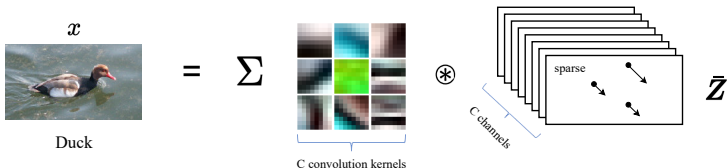
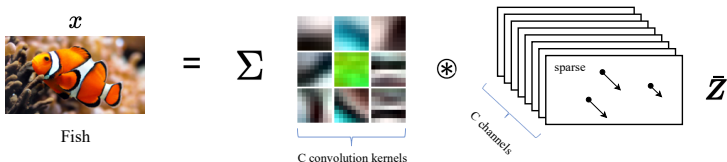
# Tradeoff between Invariance and Separability

**A problem with separability:** superposition of shifted “delta” functions can generate any other signals:  
 $\text{span}[\text{circ}(x)] = \mathbb{R}^n$ !



**A necessary assumption:**  $x$  is **sparsely generated** from incoherent dictionaries for different classes:

$$x = [\text{circ}(\mathcal{D}_1), \text{circ}(\mathcal{D}_2), \dots, \text{circ}(\mathcal{D}_k)] \bar{z}.$$

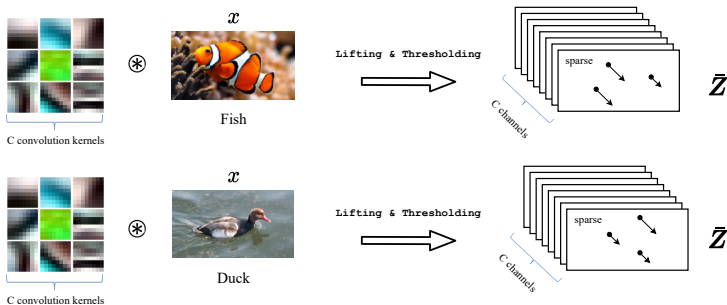


# Tradeoff between Invariance and Separability

**A basic idea:** estimate sparse codes  $\bar{z}$  by taking their responses to multiple analysis filters  $k_1, \dots, k_C \in \mathbb{R}^n$  [Rubinstein & Elad 2014]:

$$\bar{z} = \tau[k_1 \circledast x, \dots, k_C \circledast x]^* \in \mathbb{R}^{C \times n}. \quad (25)$$

for some entry-wise “sparsity-promoting” operator  $\tau(\cdot)$ .



## Multi-Channel Convolutions

Given a set of multi-channel sparse codes  $\bar{\mathbf{Z}} = [\bar{\mathbf{z}}^1, \dots, \bar{\mathbf{z}}^m]$ , construct the ReduNet from their circulant families  $\bar{\mathbf{Z}} = [\text{circ}(\bar{\mathbf{z}}^1), \dots, \text{circ}(\bar{\mathbf{z}}^m)]$ .

### Proposition (Convolution Structures of $\bar{\mathbf{E}}$ and $\bar{\mathbf{C}}^j$ )

*The linear operator in the ReduNet:*

$$\bar{\mathbf{E}} = \alpha \left( \mathbf{I} + \alpha \sum_{i=1}^m \text{circ}(\bar{\mathbf{z}}^i) \text{circ}(\bar{\mathbf{z}}^i)^* \right)^{-1} \in \mathbb{R}^{Cn \times Cn}$$

*is a block circulant matrix and represents a multi-channel convolution:*

$$\bar{\mathbf{E}}(\bar{\mathbf{z}}) = \bar{\mathbf{e}} \circledast \bar{\mathbf{z}} \in \mathbb{R}^{Cn},$$

*where  $\bar{\mathbf{e}}$  is the first slice of  $\bar{\mathbf{E}}$ . Similarly, the operators  $\bar{\mathbf{C}}^j$  associated with subsets  $\bar{\mathbf{Z}}^j$  are also multi-channel circular convolutions.*

# Multi-Channel Convolutions

$$\bar{E}(\bar{z}) = \bar{e} \circledast \bar{z} \in \mathbb{R}^{Cn}, \quad \bar{C}^j(\bar{z}) = \bar{c}^j \circledast \bar{z} \in \mathbb{R}^{Cn} :$$

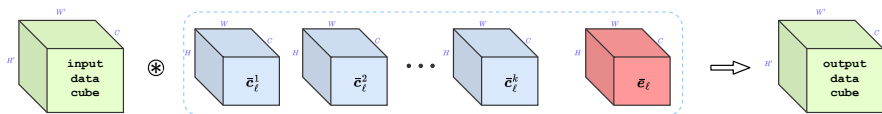


Figure:  $\bar{E}$  and  $\bar{C}^j$  are automatically multi-channel convolutions!

# The Convolution ReduNet versus Scattering Network

Iterative projected gradient ascent (PGA) for invariant rate reduction:

$$\bar{z}_{\ell+1} \propto \bar{z}_{\ell} + \eta \cdot \underbrace{\left[ \bar{E}_{\ell} \bar{z}_{\ell} + \sigma([\bar{C}_{\ell}^1 \bar{z}_{\ell}, \dots, \bar{C}_{\ell}^k \bar{z}_{\ell}]) \right]}_{g(\bar{z}_{\ell}, \theta_{\ell})}, \quad (26)$$

with each layer being a fixed number of multi-channel convolutions!

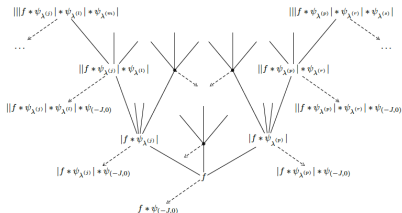
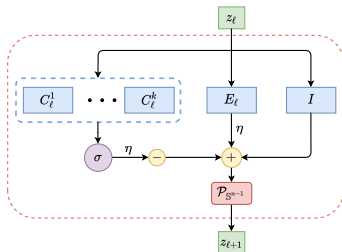


Fig. 2: Scattering network architecture based on wavelet filters and the modulus non-linearity. The elements of the feature vector  $\Phi_W(f)$  in (1) are indicated at the tips of the arrows.

Figure: Left: **ReduNet** layer. Right: **Scattering Network** [J. Bruna and S. Mallat, 2013] [T. Wiatowski and H. Blcskei, 2018] (only 2-3 layers).

# Fast Computation in Spectral Domain

**Fact:** all circulant matrices can be simultaneously diagonalized by the *discrete Fourier transform*  $\mathbf{F}$ :  $\text{circ}(\mathbf{z}) = \mathbf{F}^* \mathbf{D} \mathbf{F}$ .

$$\left( \mathbf{I} + \sum_{i=1}^m \text{circ}(\bar{\mathbf{z}}^i) \text{circ}(\bar{\mathbf{z}}^i)^* \right)^{-1} = \left( \mathbf{I} + \begin{bmatrix} \mathbf{F}^* & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{F}^* \end{bmatrix} \begin{bmatrix} \mathbf{D}_{11} & \cdots & \mathbf{D}_{1C} \\ \vdots & \ddots & \vdots \\ \mathbf{D}_{C1} & \cdots & \mathbf{D}_{CC} \end{bmatrix} \begin{bmatrix} \mathbf{F} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{F} \end{bmatrix} \right)^{-1} \in \mathbb{R}^{nC \times nC}$$

where  $\mathbf{D}_{ij}$  are all diagonal of size  $n$ .

Computing the inverse is  $O(C^3 n)$  in the spectral domain, instead of  $O(C^3 n^3)$ ! *Learning convolutional networks for invariant classification is naturally far more efficient in the spectral domain!*

**Nature:** In visual cortex, neurons encode and transmit information in frequency, hence called “spiking neurons” [Softky & Koch, 1993; Eliasmith & Anderson, 2003].

## A “White Box” Deep Convolutional ReduNet by Construction (Spectral Domain)

Require:  $\bar{Z} \in \mathbb{R}^{C \times T \times m}$ ,  $\Pi$ ,  $\epsilon > 0$ ,  $\lambda$ , and a learning rate  $\eta$ .

- 1: Set  $\alpha = \frac{C}{m\epsilon^2}$ ,  $\{\alpha_j = \frac{C}{\text{tr}(\Pi^j)\epsilon^2}\}_{j=1}^k$ ,  $\{\gamma_j = \frac{\text{tr}(\Pi^j)}{m}\}_{j=1}^k$ .
- 2: Set  $\bar{V}_0 = \{\bar{v}_0^i(p) \in \mathbb{C}^C\}_{p=0, i=1}^{T-1, m} \doteq \text{DFT}(\bar{Z}) \in \mathbb{C}^{C \times T \times m}$ .
- 3: **for**  $\ell = 1, 2, \dots, L$  **do**
- 4:   **for**  $p = 0, 1, \dots, T-1$  **do**
- 5:     Compute  $\bar{\mathcal{E}}_\ell(p) \in \mathbb{C}^{C \times C}$  and  $\{\bar{\mathcal{C}}_\ell^j(p) \in \mathbb{C}^{C \times C}\}_{j=1}^k$  as  

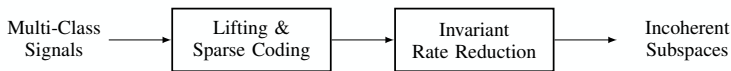
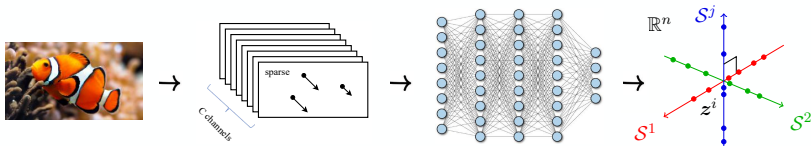
$$\bar{\mathcal{E}}_\ell(p) \doteq \alpha \cdot [\mathbf{I} + \alpha \cdot \bar{V}_{\ell-1}(p) \cdot \bar{V}_{\ell-1}(p)^*]^{-1},$$

$$\bar{\mathcal{C}}_\ell^j(p) \doteq \alpha_j \cdot [\mathbf{I} + \alpha_j \cdot \bar{V}_{\ell-1}(p) \cdot \Pi^j \cdot \bar{V}_{\ell-1}(p)^*]^{-1};$$
- 6:   **end for**
- 7:   **for**  $i = 1, \dots, m$  **do**
- 8:     **for**  $p = 0, 1, \dots, T-1$  **do**
- 9:       Compute  $\{\bar{\mathcal{P}}_\ell^{ij}(p) \doteq \bar{\mathcal{C}}_\ell^j(p) \cdot \bar{v}_\ell^i(p) \in \mathbb{C}^{C \times 1}\}_{j=1}^k$ ;
- 10:     **end for**
- 11:     Let  $\{\bar{\mathcal{P}}_\ell^{ij} = [\bar{\mathcal{P}}_\ell^{ij}(0), \dots, \bar{\mathcal{P}}_\ell^{ij}(T-1)] \in \mathbb{C}^{C \times T}\}_{j=1}^k$ ;
- 12:     Compute  $\{\hat{\pi}_\ell^{ij} = \frac{\exp(-\lambda \|\bar{\mathcal{P}}_\ell^{ij}\|_F)}{\sum_{j=1}^k \exp(-\lambda \|\bar{\mathcal{P}}_\ell^{ij}\|_F)}\}_{j=1}^k$ ;
- 13:     **for**  $p = 0, 1, \dots, T-1$  **do**
- 14:       
$$\bar{v}_\ell^i(p) = \bar{v}_{\ell-1}^i(p) + \eta \left( \bar{\mathcal{E}}_\ell(p) \bar{v}_\ell^i(p) - \sum_{j=1}^k \gamma_j \cdot \hat{\pi}_\ell^{ij} \cdot \bar{\mathcal{C}}_\ell^j(p) \cdot \bar{v}_\ell^i(p) \right);$$
- 15:     **end for**
- 16:     
$$\bar{v}_\ell^i = \bar{v}_\ell^i / \|\bar{v}_\ell^i\|_F;$$
- 17:   **end for**
- 18:   Set  $\bar{Z}_\ell = \text{IDFT}(\bar{V}_\ell)$  as the feature at the  $\ell$ -th layer;
- 19:   
$$\frac{1}{2T} \sum_{p=0}^{T-1} \left( \log \det[\mathbf{I} + \alpha \bar{V}_\ell(p) \cdot \bar{V}_\ell(p)^*] - \frac{\text{tr}(\Pi^j)}{m} \log \det[\mathbf{I} + \alpha_j \bar{V}_\ell(p) \cdot \Pi^j \cdot \bar{V}_\ell(p)^*] \right);$$
- 20: **end for**

Ensure: features  $\bar{Z}_L$ , the learned filters  $\{\bar{\mathcal{E}}_\ell(p)\}_{\ell, p}$  and  $\{\bar{\mathcal{C}}_\ell^j(p)\}_{j, \ell, p}$ .

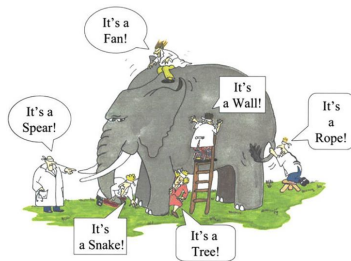


# Overall Process (the Elephant)



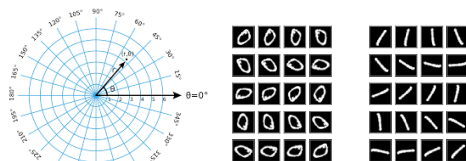
## Necessary components:

- **sparse coding** for class separability;
- **deep networks** maximize rate reduction;
- **spectral computing** for shift-invariance;
- convolution, normalization, nonlinearity...

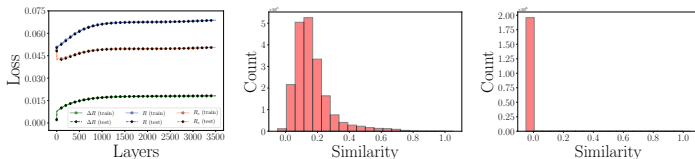


# Experiment: 1D Cyclic Shift Invariance of 0 and 1

2000 training samples, 1980 testing,  $C = 5$ ,  $L = 3500$ -layers ReduNet.

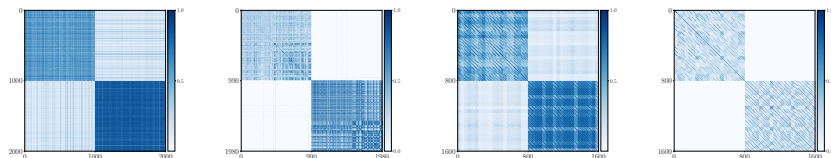


**Figure:** Left: Multi-channel feature representation of an image in polar coordinates. Right: Example of training/testing samples.



**Figure:** Left: Rates along the layers; Middle: cross-class cosine similarity among trainings; Right: similarity among testings.

# Experiment: 1D Cyclic Shift Invariance of 0 and 1



**Figure:** Left two: heat maps for training and testing. Right two: heat maps for one pair of samples at every possible shift.

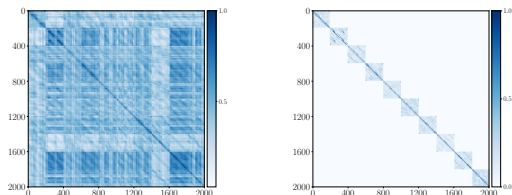
**Table:** Network performance on digits with **all rotations**.

	REDUNET	REDUNET (INVARIANT)
ACC (ORIGINAL TEST DATA)	0.983	0.996
ACC (TEST WITH ALL SHIFTS)	0.707	0.993

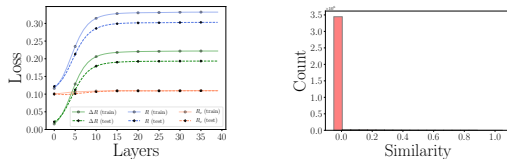
1. Fooling CNNs with simple transformations, Engstrom et.al., 2017.
2. Why do deep convolutional networks generalize so poorly to small image transformations? Azulay & Weiss, 2018.

# Experiment: 1D Cyclic Shift Invariance of All 10 Digits

100 training samples, 100 testing,  $C = 20$ ,  $L = 40$ -layers ReduNet.



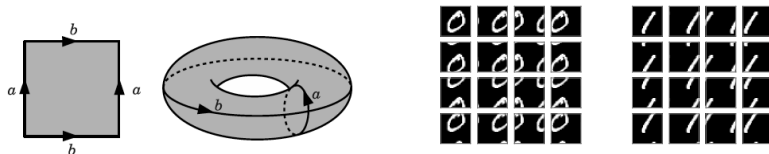
**Figure:** Heatmaps of cosine similarity among shifted training data  $\mathbf{X}_{\text{shift}}$  (left) and learned features  $\mathbf{Z}_{\text{shift}}$  (right).



**Figure:** Left: Rates evolution with iterations; Right: histograms of the cosine similarity (in absolute value) between all pairs of features across different classes.

# Experiment: 2D Cyclic Translation Invariance

1000 for training, 500 for testing,  $C = 5$ ,  $L = \mathbf{2000}$ -layers ReduNet.



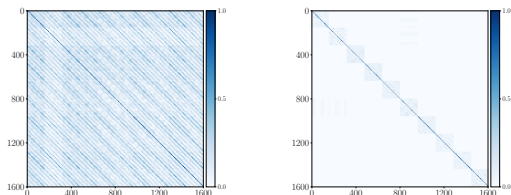
**Table:** Network performance on digits with **all translations**.

	REDUNET	REDUNET (INVARIANT)
ACC (ORIGINAL TEST DATA)	0.980	0.975
ACC (TEST WITH ALL SHIFTS)	0.540	0.909

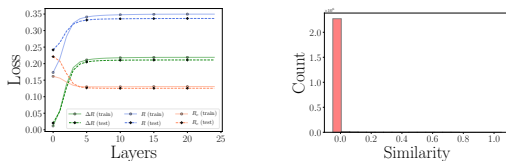
1. Fooling CNNs with simple transformations, Engstrom et.al., 2017.
2. Why do deep convolutional networks generalize so poorly to small image transformations? Azulay & Weiss, 2018.

# Experiment: 2D Cyclic Trans. Invariance of All 10 Digits

100 training samples, 100 testing,  $C = 75$ ,  $L = 25$ -layers ReduNet.



**Figure:** Heatmaps of cosine similarity among shifted training data  $\mathbf{X}_{\text{shift}}$  (left) and learned features  $\mathbf{Z}_{\text{shift}}$  (right).



**Figure:** Left: Rates evolution with iterations; Right: histograms of the cosine similarity (in absolute value) between all pairs of features across different classes.

## Experiment: Back Propagation of ReduNet (MNIST)

2D cyclic trans. of 10 digits, 500 training samples, all testing,  $C = 16$ ,  $L = 30$ -layers invariant ReduNet.

Initialization	Backpropagation	Test Accuracy
✓	✗	89.8%
✗	✓	93.2%
✓	✓	97.8%

**Table:** Test accuracy of 2D translation-invariant ReduNet, ReduNet-bp (without initialization), and ReduNet-bp (with initialization) on the MNIST dataset.

- **Backprop:** the ReduNet architecture *can* be fine-tuned by SGD and achieves better standard accuracy after back propagation;
- **Initialization:** using ReduNet for initialization can achieve better performance than the same architecture with random initialization.

# Experiment: Back Propagation of ReduNet (CIFAR-10)

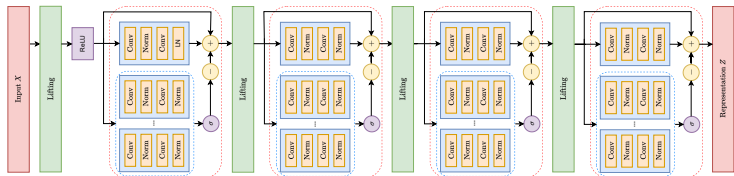


Figure: A ReduNet-inspired architecture

Table: Classification performance of ReduNet-inspired architecture on CIFAR10.

ReLU	TRAIN ACC	TEST ACC
✓	0.9997	0.8327
✗	0.9970	0.6542



# Recap: White-Box Deep Networks

**A promising approach:** signal models  $\Rightarrow$  deep architectures

- Scattering networks [Bruna & Mallat 2013]
- Convolutional sparse coding networks [Papayan et al. 2018]
- ReduNets [Chan, Yu et al. 2022]

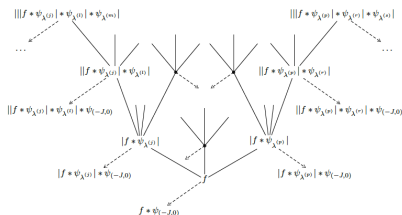
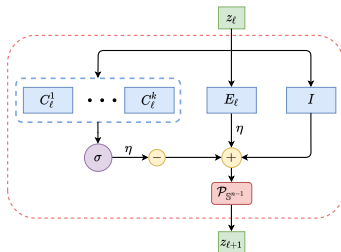


Fig. 2: Scattering network architecture based on wavelet filters and the modulus non-linearity. The elements of the feature vector  $\Phi_W(f)$  in (1) are indicated at the tips of the arrows.

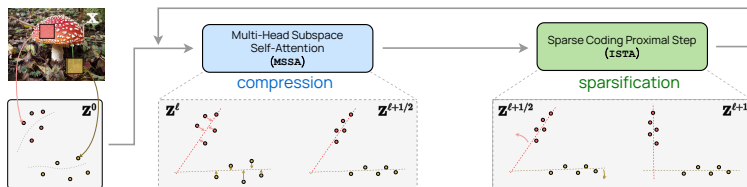
Figure: Left: **ReduNet** layer. Right: **Scattering Network** [Bruna & Mallat 2013] [Wiatowski & Bölcskei 2018] (only 2-3 layers).

**Pitfall of existing methods: Challenging to scale to massive datasets with strong performance**

# CRATE: A White-Box Transformer via Sparse MCR<sup>2</sup>

A **white-box**, **mathematically interpretable**, **transformer-like** deep network architecture from **iterative unrolling** optimization schemes to incrementally optimize the sparse rate reduction objective:

$$\max_{f \in \mathcal{F}} \mathbb{E}_{\mathbf{Z}} [\Delta R(\mathbf{Z}; \mathbf{U}_{[K]}) - \|\mathbf{Z}\|_0], \quad \mathbf{Z} = f(\mathbf{X}).$$



**CRATE: White-Box Transformers via Sparse Rate Reduction**

<https://arxiv.org/abs/2306.01129>

# Conclusions: Learn to Compress and Compress to Learn!

## Principles of Parsimony:

- Clustering via compression:  $\min_{\Pi} R^c(\mathbf{X}, \Pi)$
- Classification via compression:  $\min_{\pi} \delta R^c(\mathbf{x}, \pi)$
- Representation via maximizing rate reduction:  $\max_{\mathbf{Z}} \Delta R(\mathbf{Z}, \Pi)$
- Deep networks via optimizing rate reduction:  $\dot{\mathbf{Z}} = \eta \cdot \frac{\partial \Delta R}{\partial \mathbf{Z}}$

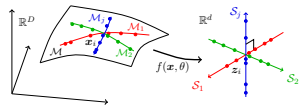
## A Unified Framework:

- A principled objective for all settings of learning: **information gain**
- A principled approach to interpret deep networks: **optimization**

*“Everything should be made as simple as possible, but not simpler.”*  
– Albert Einstein

# White-box Objectives, Architectures, and Representations

Comparison with conventional practice of NNs (since McCulloch-Pitts'1943).



	Conventional DNNs	ReduNets
Objectives	input/output fitting	information gain
Deep architectures	trial & error	iterative optimization
Layer operators	empirical	projected gradient
Shift invariance	CNNs+augmentation	invariant ReduNets
Initializations	random/pre-design	forward unrolled <sup>3</sup>
Training/fine-tuning	back prop	forward/back prop
Interpretability	black box	white box
Representations	hidden/latent	incoherent subspaces

<sup>3</sup> *The Forward-Forward Algorithm*, G. Hinton, 2022.

# Open Problems: Theory

$$\mathbf{MCR}^2: \max_{\mathbf{Z} \in \mathbb{S}^{d-1}, \mathbf{\Pi} \in \Omega} \Delta R(\mathbf{Z}, \mathbf{\Pi}, \epsilon) = R(\mathbf{Z}, \epsilon) - R^c(\mathbf{Z}, \epsilon \mid \mathbf{\Pi}).$$

- **Phase transition** phenomenon in clustering via compression?
- Statistical justification for **robustness** of  $\mathbf{MCR}^2$  to label noise?
- **Optimal configurations** for broader conditions and distributions?
- Fundamental **tradeoff** between sparsity and invariance?
- **Jointly optimizing** both representation  $\mathbf{Z}$  and clustering  $\mathbf{\Pi}$ ?

$$\mathbf{Joint Dynamics:} \quad \dot{\mathbf{Z}} = \eta \cdot \frac{\partial \Delta R}{\partial \mathbf{Z}}, \quad \dot{\mathbf{\Pi}} = \gamma \cdot \frac{\partial \Delta R}{\partial \mathbf{\Pi}}.$$

# Open Problems: Architectures and Algorithms

## Gradient of Rate Distortion:

$$\left. \frac{\partial R(\mathbf{Z})}{\partial \mathbf{Z}} \right|_{\mathbf{Z}_\ell} = \underbrace{\alpha(\mathbf{I} + \alpha \mathbf{Z}_\ell \mathbf{Z}_\ell^*)^{-1} \mathbf{Z}_\ell}_{\text{auto-regress residual}} \approx \underbrace{\alpha[\mathbf{Z}_\ell - \alpha \mathbf{Z}_\ell (\mathbf{Z}_\ell^* \mathbf{Z}_\ell)]}_{\text{self-attention head}}.$$

## ReduNet:

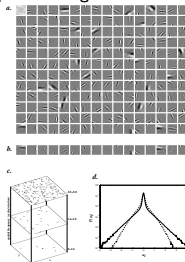
$$\bar{\mathbf{z}}_{\ell+1} \propto \bar{\mathbf{z}}_\ell + \eta \cdot \left[ \bar{\mathbf{e}}_\ell \circledast \bar{\mathbf{z}}_\ell + \sigma([\bar{\mathbf{c}}_\ell^1 \circledast \bar{\mathbf{z}}_\ell, \dots, \bar{\mathbf{c}}_\ell^k \circledast \bar{\mathbf{z}}_\ell]) \right] \in \mathbb{S}^{d-1}.$$

- New architectures from **accelerated** gradient schemes?
- Conditions for channel-wise **separable and short** convolutions?
- Architectures from invariant rate reduction for **other groups**?
- Algorithmic architectures (or networks) for optimizing  $\mathbf{Z}, \mathbf{\Pi}$  jointly?

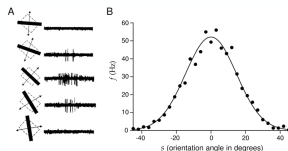
# Open Directions: Extensions

- Data with other **dynamical or graphical** structures.
- Better transferability and robustness w.r.t. **low-dim structures**.
- Combine with a **generative model** (a generator or decoder).
- Sparse coding, spectral computing, subspace embedding in **nature**.<sup>4</sup>

## sparse coding in visual cortex



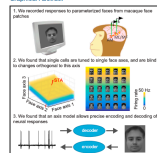
**Rate coding hypothesis:** the signal conveyed by a neuron is in the *rate* of spiking. Spiking irregularity is largely due to noise and does not convey information.



## Cell

### The Code for Facial Identity in the Primate Brain

#### Graphical Abstract



#### Authors

Le Chang, Doris Y. Tsao

Correspondence  
lechang@caltech.edu (L.C.),  
dtsao@caltech.edu (D.Y.T.)

#### In Brief

Facial identity is encoded via a remarkably simple neural code that relies on the ability of neurons to distinguish facial features along specific axes in face space, preserving the long-standing assumption that single face cells encode individual faces.

#### Highlights

- Facial images can be linearly reconstructed using responses of ~200 face cells.
- Face cells display flat tuning along dimensions orthogonal to the axis being coded.
- The axis model is more efficient, robust, and flexible than the exemplar model.
- Face patches MIPF and AM carry complementary information about faces.

<sup>4</sup>figures from Bruno Olshausen of Neuroscience Dept., UC Berkeley.

# References: White-Box Deep Networks via Rate Reduction

- ① **ReduNet**: A Whitebox Deep Network from Rate Reduction (JMLR 2022):  
<https://arxiv.org/abs/2105.10446>
- ② **Representation** via Maximal Coding Rate Reduction (NeurIPS 2020):  
<https://arxiv.org/abs/2006.08558>
- ③ **Classification** via Minimal Incremental Coding Length (NIPS 2007):  
[http://people.eecs.berkeley.edu/~yima/psfile/MICL\\_SJIS.pdf](http://people.eecs.berkeley.edu/~yima/psfile/MICL_SJIS.pdf)
- ④ **Clustering** via Lossy Coding and Compression (TPAMI 2007):  
<http://people.eecs.berkeley.edu/~yima/psfile/Ma-PAMI07.pdf>



# Source Code: Whitebox ReduNet

① **Github Link:**

<https://github.com/Ma-Lab-Berkeley/ReduNet>

② **Google Colab:**

[https://colab.research.google.com/github/ryanchankh/redunet\\_demo/blob/master/gaussian3d.ipynb](https://colab.research.google.com/github/ryanchankh/redunet_demo/blob/master/gaussian3d.ipynb)

③ **Jupyter Notebook:**

[https://github.com/ryanchankh/redunet\\_demo/blob/master/gaussian3d.ipynb](https://github.com/ryanchankh/redunet_demo/blob/master/gaussian3d.ipynb)

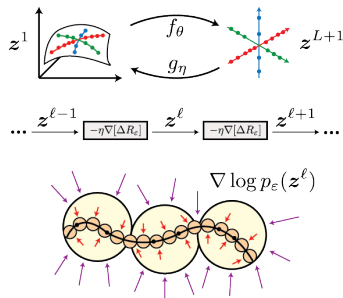
## Textbook

## Learning Deep Representations of Data Distributions

Sam Buchanan, Druv Pai, Peng Wang, and Yi Ma  
Version 1.0, August 18, 2025.

An open-source book on the GitHub:

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



# Deep (Convolution) Network Architectures are Iterative Optimization for Compression!

Thank you!  
Questions, please?

*"What I cannot create, I do not understand."*  
– Richard Feynman



SIMONS  
FOUNDATION