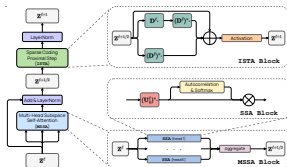


# White-Box Transformers via Sparse Rate Reduction

Yaodong Yu, Sam Buchanan, Druv Pai,  
Tianzhe Chu, Ziyang Wu, Shengbang Tong,  
Benjamin D. Haeffele, Yi Ma

June 10, 2023



# Outline

## ① CRATE

White-Box Architectures for Representation Learning  
CRATE: White-Box Transformers from Sparse MCR<sup>2</sup>  
Experimental Results on CRATE

## ② Conclusion

# Identification/Representation of High-Dim Structured Data

*Focus on one half of our goal:*

Given samples

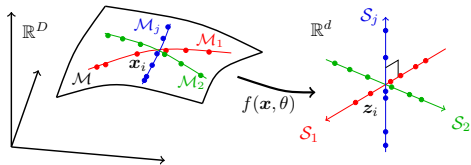
$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m] \subset \cup_{j=1}^k \mathcal{M}_j,$$

seek a good representation

$$\mathbf{Z} = [z_1, \dots, z_m] \subset \mathbb{R}^d$$

through a continuous mapping:

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



**How to obtain a white-box architecture  $f$  that simultaneously identifies and represents large-scale datasets?**

## Recap: White-Box Deep Networks

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

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

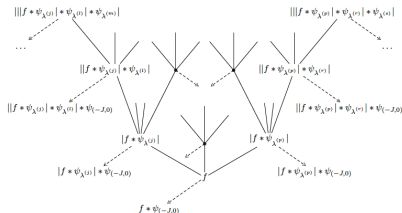
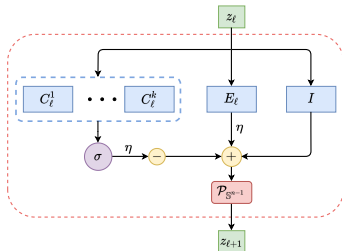


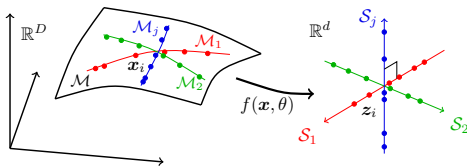
Fig. 2. Scattering network architecture based on wavelet filters and the modulus non-linearity. The elements of the feature vector  $\Phi_W(f)$  in (1) are indicated at the tips of the arrows.

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

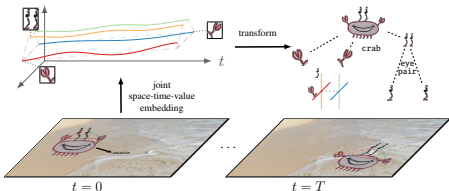
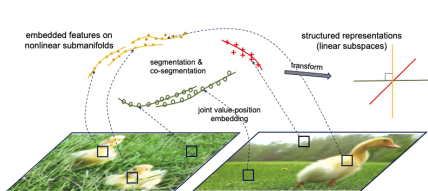
**Pitfall of existing methods: Challenging to scale to massive datasets with strong performance**

# Improved White-Box Scaling by Improved Signal Modeling?

So far: *Each sample is drawn from a mixture of manifolds*



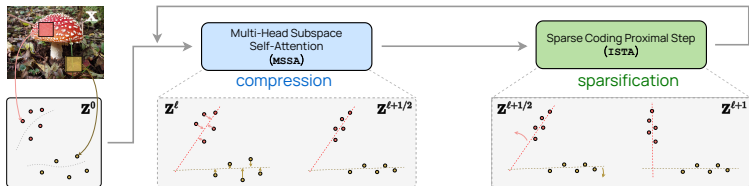
Better? *Each sample  $\supset$  correlated tokens—mixture of manifold marginals!*



# 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**

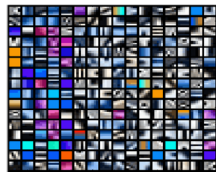
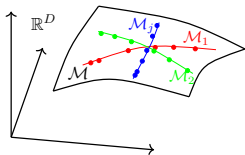
<https://arxiv.org/abs/2306.01129>

# Sparse MCR<sup>2</sup> Objective and Incremental Representation

The sparse rate reduction (**Sparse MCR<sup>2</sup>**) objective is defined as

$$\begin{aligned} & \arg \max_{f \in \mathcal{F}} \mathbb{E}_{\mathbf{Z}} [\Delta R(\mathbf{Z}; \mathbf{U}_{[K]}) - \|\mathbf{Z}\|_0] \\ &= \arg \min_{f \in \mathcal{F}} \mathbb{E}_{\mathbf{Z}} \left[ \underbrace{R^c(\mathbf{Z}; \mathbf{U}_{[K]})}_{\text{compression}} + \underbrace{\|\mathbf{Z}\|_0 - R(\mathbf{Z})}_{\text{sparsification}} \right]. \end{aligned}$$

$\mathbf{U}_{[K]} = (\mathbf{U}_1, \dots, \mathbf{U}_K)$ ,  $\mathbf{U}_k \in \mathbb{R}^{d \times p}$  are *subspaces parameterizing the marginal distribution of tokens*  $(\mathbf{z}_i)_{i=1}^N$



# Sparse MCR<sup>2</sup> Objective and Incremental Representation

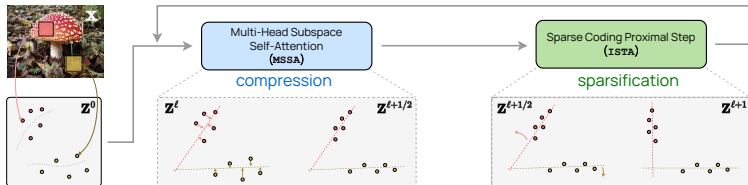
The sparse rate reduction (**Sparse MCR<sup>2</sup>**) objective is defined as

$$\begin{aligned} & \arg \max_{f \in \mathcal{F}} \mathbb{E}_{\mathbf{Z}} [\Delta R(\mathbf{Z}; \mathbf{U}_{[K]}) - \|\mathbf{Z}\|_0] \\ &= \arg \min_{f \in \mathcal{F}} \mathbb{E}_{\mathbf{Z}} \left[ \underbrace{R^c(\mathbf{Z}; \mathbf{U}_{[K]})}_{\text{compression}} + \underbrace{\|\mathbf{Z}\|_0 - R(\mathbf{Z})}_{\text{sparsification}} \right]. \end{aligned}$$

The global transformation  $f$  is realized through **local transformations**:

$$f: \mathbf{X} \xrightarrow{f^0} \mathbf{Z}^0 \rightarrow \dots \rightarrow \mathbf{Z}^\ell \xrightarrow{f^\ell} \mathbf{Z}^{\ell+1} \rightarrow \dots \rightarrow \mathbf{Z}^L = \mathbf{Z}.$$

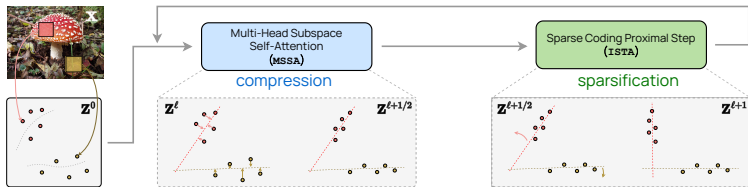
Each  $f^\ell$  deforms  $\mathbf{Z}^\ell$  according to its own **local signal model**  $\mathbf{U}_{[K]}^\ell$ .



# Sparse MCR<sup>2</sup> Objective and Incremental Representation

The sparse rate reduction (Sparse MCR<sup>2</sup>) objective is defined as

$$\begin{aligned} & \arg \max_{f \in \mathcal{F}} \mathbb{E}_{\mathbf{Z}} [\Delta R(\mathbf{Z}; \mathbf{U}_{[K]}) - \|\mathbf{Z}\|_0] \\ &= \arg \min_{f \in \mathcal{F}} \mathbb{E}_{\mathbf{Z}} \left[ \underbrace{R^c(\mathbf{Z}; \mathbf{U}_{[K]})}_{\text{compression}} + \underbrace{\|\mathbf{Z}\|_0 - R(\mathbf{Z})}_{\text{sparsification}} \right]. \end{aligned}$$



How to construct a representation  $f$  to incrementally optimize the **compression** term and the **sparsification** term?

# Recap: Compression and Expansion in MCR<sup>2</sup>

## Expansion:

$$R(\mathbf{Z}) = \frac{1}{2} \sum_{k=1}^K \log \det \left( \mathbf{I} + \frac{d}{N\epsilon^2} \mathbf{Z}^* \mathbf{Z} \right)$$

## Compression:

$$R^c(\mathbf{Z}; \mathbf{U}_{[K]}) = \frac{1}{2} \sum_{k=1}^K \log \det \left( \mathbf{I} + \frac{p}{N\epsilon^2} (\mathbf{U}_k^* \mathbf{Z})^* (\mathbf{U}_k^* \mathbf{Z}) \right)$$

## Gradient of Rate Distortion:

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

## Compression in Sparse MCR<sup>2</sup>

To optimize the compression term  $R^c(\mathbf{Z}; \mathbf{U}_{[K]})$ , we propose to compress the set of tokens against the subspaces  $(\mathbf{U}_k)_{k=1}^K$  by minimizing the coding rate via “approximate” gradient descent

$$\begin{aligned} \text{(Gradient Descent): } \quad \mathbf{Z}^\ell - \kappa \nabla_{\mathbf{Z}} R^c(\mathbf{Z}^\ell; \mathbf{U}_{[K]}) \\ \approx \left(1 - \kappa \cdot \frac{p}{N\epsilon^2}\right) \mathbf{Z}^\ell + \kappa \cdot \frac{p}{N\epsilon^2} \cdot \text{MSSA}(\mathbf{Z}^\ell | \mathbf{U}_{[K]}), \end{aligned}$$

where MSSA is defined through an SSA operator as:

$$\begin{aligned} \text{SSA}(\mathbf{Z} | \mathbf{U}_k) &= (\mathbf{U}_k^* \mathbf{Z}) \text{softmax}((\mathbf{U}_k^* \mathbf{Z})^* (\mathbf{U}_k^* \mathbf{Z})), \\ \text{MSSA}(\mathbf{Z} | \mathbf{U}_{[K]}) &= \frac{p}{N\epsilon^2} \cdot [\mathbf{U}_1, \dots, \mathbf{U}_K] \begin{bmatrix} \text{SSA}(\mathbf{Z} | \mathbf{U}_1) \\ \vdots \\ \text{SSA}(\mathbf{Z} | \mathbf{U}_K) \end{bmatrix}. \end{aligned}$$

**No need for separate query- $Q$ , key- $K$ , value- $V$  in transformer attention block.**

## Compression in Sparse MCR<sup>2</sup>

To optimize the compression term  $R^c(\mathbf{Z}; \mathbf{U}_{[K]})$ , we propose to compress the set of tokens against the subspaces  $(\mathbf{U}_k)_{k=1}^K$  by minimizing the coding rate via “approximate” gradient descent

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

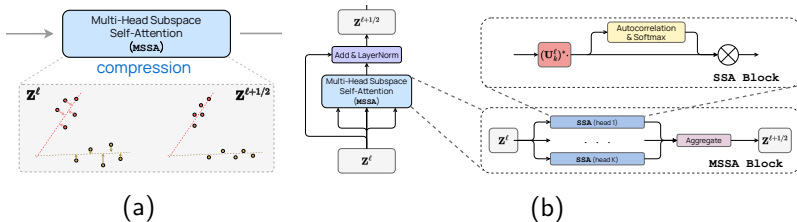


Figure: (a). Visualization of MSSA block; (b). Architecture of MSSA block.

## Sparsification in Sparse MCR<sup>2</sup>

To optimize the sparsification term  $\|\mathbf{Z}\|_0 - R(\mathbf{Z})$ , we posit a incoherent or orthogonal dictionary  $\mathbf{D} \in \mathbb{R}^{d \times d}$  and sparsify  $\mathbf{Z}^{\ell+1/2}$  with respect to  $\mathbf{D}$ , that is

$$\mathbf{Z}^{\ell+1/2} = \mathbf{D}\mathbf{Z}^{\ell+1}.$$

By the incoherence assumption, we have  $\mathbf{D}^*\mathbf{D} \approx \mathbf{I}_d$ ; thus

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

Thus we approximately optimize the **sparsification objective** with the following program:

$$\mathbf{Z}^{\ell+1} = \operatorname{argmin}_{\mathbf{Z}} \|\mathbf{Z}\|_0 \quad \text{subject to} \quad \mathbf{Z}^{\ell+1/2} = \mathbf{D}\mathbf{Z}.$$

## Sparsification in Sparse MCR<sup>2</sup>

Given the sparse representation program

$$\mathbf{Z}^{\ell+1} = \operatorname{argmin}_{\mathbf{Z}} \|\mathbf{Z}\|_0 \quad \text{subject to} \quad \mathbf{Z}^{\ell+1/2} = \mathbf{D}\mathbf{Z}.$$

we can relax it to an convex program, i.e., **positive sparse coding**:

$$\mathbf{Z}^{\ell+1} = \operatorname{argmin}_{\mathbf{Z} \geq 0} \left[ \lambda \|\mathbf{Z}\|_1 + \|\mathbf{Z}^{\ell+1/2} - \mathbf{D}\mathbf{Z}\|_F^2 \right].$$

We can incrementally optimize the above objective by performing an unrolled proximal gradient descent step, known as an ISTA step:

$$\begin{aligned} \mathbf{Z}^{\ell+1} &= \operatorname{ReLU}(\mathbf{Z}^{\ell+1/2} + \eta \mathbf{D}^* (\mathbf{Z}^{\ell+1/2} - \mathbf{D}\mathbf{Z}^{\ell+1/2}) - \eta \lambda \mathbf{1}) \\ &:= \operatorname{ISTA}(\mathbf{Z}^{\ell+1/2} \mid \mathbf{D}^{\ell}). \end{aligned}$$

**The ISTA block uses much fewer parameters than transformer MLP block, and provides more interpretable representations.**

## Sparsification in Sparse MCR<sup>2</sup>

To optimize the sparsification term  $\|\mathbf{Z}\|_0 - R(\mathbf{Z})$ , we propose to apply an unrolled proximal gradient descent step, known as an ISTA step:

$$\begin{aligned}\mathbf{Z}^{\ell+1} &= \text{ReLU}(\mathbf{Z}^{\ell+1/2} + \eta \mathbf{D}^* (\mathbf{Z}^{\ell+1/2} - \mathbf{D} \mathbf{Z}^{\ell+1/2}) - \eta \lambda \mathbf{1}) \\ &:= \text{ISTA}(\mathbf{Z}^{\ell+1/2} \mid \mathbf{D}^\ell).\end{aligned}$$

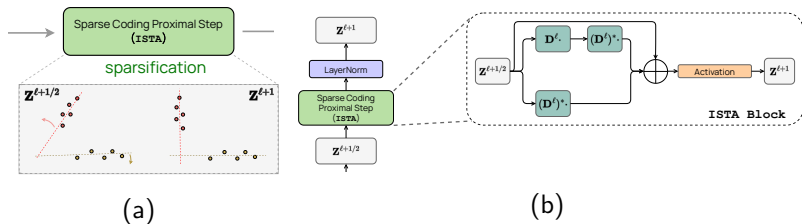


Figure: (a). Visualization of ISTA block; (b). Architecture of ISTA block.

# One Layer of CRATE

Each layer of **CRATE** thus incrementally optimizes the **compression term**  $R^c(\mathbf{Z}; \mathbf{U}_{[K]})$  and **sparsification term**  $\|\mathbf{Z}\|_0 - R(\mathbf{Z})$ ,

$$\mathbf{Z}^{\ell+1} = f^\ell(\mathbf{Z}^\ell) = \text{ISTA}\left(\underbrace{(\text{Id} + \text{MSSA})(\mathbf{Z}^\ell)}_{\mathbf{Z}^{\ell+1/2}}\right).$$

More specifically,

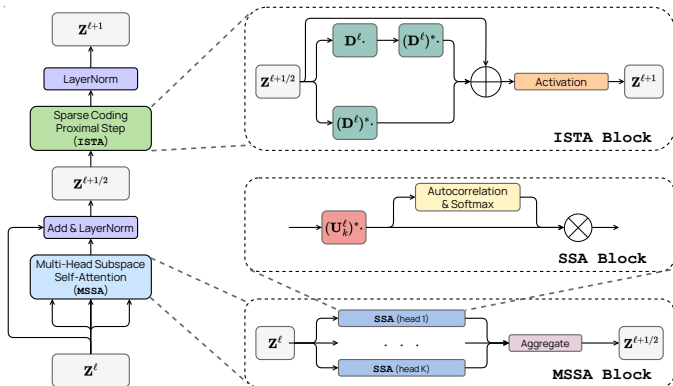
$$\mathbf{Z}^{\ell+1/2} = \mathbf{Z}^\ell + \text{MSSA}(\mathbf{Z}^\ell | \mathbf{U}_{[K]}^\ell), \quad [\text{Compression step}]$$

$$\mathbf{Z}^{\ell+1} = \text{ISTA}(\mathbf{Z}^{\ell+1/2} | \mathbf{D}^\ell), \quad [\text{Sparsification step}]$$

so the  $\ell$ -th layer of the global representation  $f$  is

$$f^\ell: \mathbf{Z}^\ell \xrightarrow{\text{Id+MSSA}} \mathbf{Z}^{\ell+1/2} \xrightarrow{\text{ISTA}} \mathbf{Z}^{\ell+1}.$$

# Overall White-Box CRATE Architecture



- Forward optimization: perform **compression** and **sparsification**.
- Learning from data: apply SGD to learn  $(\mathbf{U}_{[K]}^\ell, \mathbf{D}^\ell)_{\ell=1}^L$  from data.

# Experiment I: Supervised Learning on ImageNet-1K

**Experimental setup:** let the CLS token of  $Z^L$  (i.e., the output token set of the last layer), and then apply a linear linear to perform supervised learning on ImageNet-1K using our proposed CRATE architecture.

**Table 1:** Top 1 accuracy of CRATE on various datasets with different model scales when pre-trained on ImageNet. For ImageNet/ImageNetReal, we directly evaluate the top-1 accuracy. For other datasets, we use models that are pre-trained on ImageNet as initialization and the evaluate the transfer learning performance via fine-tuning.

Datasets	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	66.7	69.2	70.8	71.3	71.5	72.4
ImageNet Real	74.0	76.0	76.5	77.4	78.3	78.4

- CRATE demonstrates promising performance on the ImageNet-1K dataset, indicating its potential for further advancement.

# Experiment I: Supervised Learning on ImageNet-1K

**Experimental setup:** apply the CRATE model pre-trained on ImageNet-1K as initialization, and then evaluate transfer learning performance via fine-tuning.

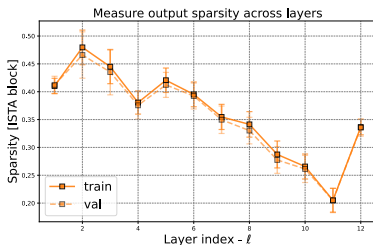
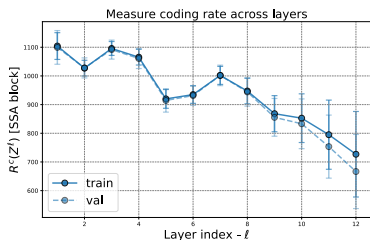
**Table 1:** Top 1 accuracy of CRATE on various datasets with different model scales when pre-trained on ImageNet. For ImageNet/ImageNetReal, we directly evaluate the top-1 accuracy. For other datasets, we use models that are pre-trained on ImageNet as initialization and the evaluate the transfer learning performance via fine-tuning.

Datasets	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	66.7	69.2	70.8	71.3	71.5	72.4
ImageNet 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

- CRATE achieves performance close to thoroughly engineered vision transformers.
- Promising scaling behavior in CRATE.

## Experiment II: Layer-wise Analysis of CRATE

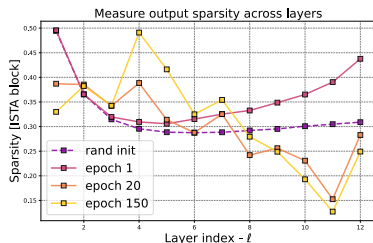
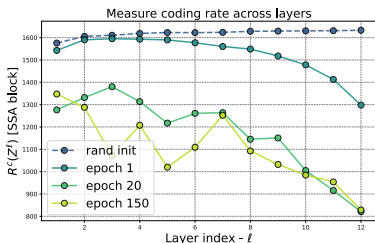
Given a learned CRATE model, we measure the compression term of  $\mathbf{Z}^{\ell+1/2}$  (left,  $R^c(\mathbf{Z}^{\ell+1/2})$ ) and the sparsification term of  $\mathbf{Z}^{\ell+1}$  (right,  $\|\mathbf{Z}^{\ell+1}\|_0$ ) on train/validation samples at **each layer**.



- The learned CRATE model indeed performs its design objective – each layer incrementally optimizes the compression term and the sparsification term.

## Experiment II: Layer-wise Analysis of CRATE

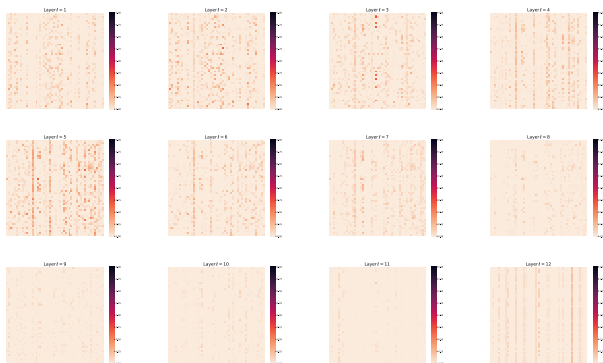
For comparison, we measure the compression/sparsification term of randomly initialized CRATE model and models at different epochs.



- Without learning from data, the random initialized CRATE model does not perform its design objective effectively.

# Experiment III: Visualize Layer-wise Output of CRATE

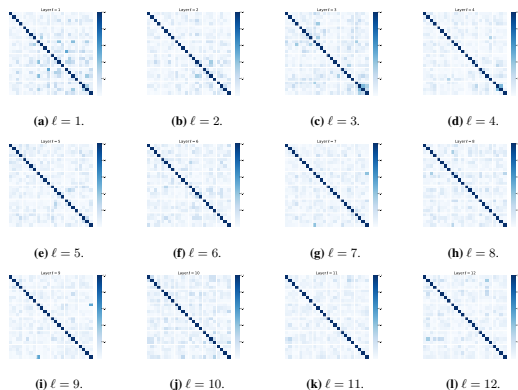
We use heatmaps to visualize the output of each layer in CRATE ( $Z^{\ell+1}$ ).



- We observe clear sparse and low-rank patterns of intermediate outputs of CRATE.

## Experiment IV: Visualize Learned Subspaces of CRATE

We use heatmaps to visualize the correlations between different subspaces  $(\mathbf{U}_k)_{k=1}^K$  of each MSSA layer in CRATE, i.e.,  $[\mathbf{U}_1^\ell, \dots, \mathbf{U}_K^\ell]^* [\mathbf{U}_1^\ell, \dots, \mathbf{U}_K^\ell]$ .



- The learned subspaces in MSSA blocks are incoherent.

# Outline

## ① CRATE

White-Box Architectures for Representation Learning  
CRATE: White-Box Transformers from Sparse MCR<sup>2</sup>  
Experimental Results on CRATE

## ② Conclusion

**Thank You! Questions?**