

---

# White-Box Transformers via Sparse Rate Reduction

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

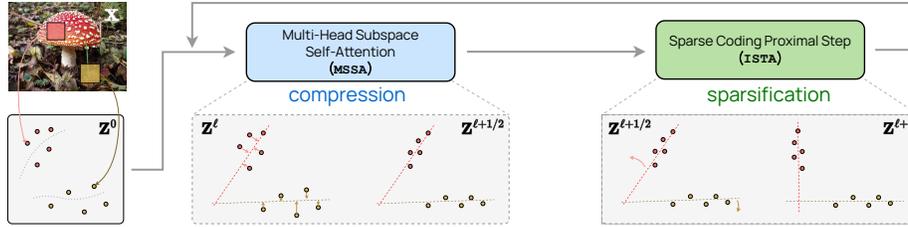
1 In this paper, we contend that the objective of representation learning is to compress  
2 and transform the distribution of the data, say sets of tokens, towards a mixture of  
3 low-dimensional Gaussian distributions supported on incoherent subspaces. The  
4 quality of the final representation can be measured by a unified objective function  
5 called *sparse rate reduction*. From this perspective, popular deep networks such  
6 as transformers can be naturally viewed as realizing iterative schemes to optimize  
7 this objective incrementally. Particularly, we show that the standard transformer  
8 block can be derived from alternating optimization on complementary parts of  
9 this objective: the multi-head self-attention operator can be viewed as a gradient  
10 descent step to compress the token sets by minimizing their lossy coding rate, and  
11 the subsequent multi-layer perceptron can be viewed as attempting to sparsify the  
12 representation of the tokens. This leads to a family of *white-box* transformer-like  
13 deep network architectures which are mathematically fully interpretable. Despite  
14 their simplicity, experiments show that these networks indeed learn to optimize  
15 the designed objective: they compress and sparsify representations of large-scale  
16 real-world vision datasets such as ImageNet, and achieve performance very close  
17 to thoroughly engineered transformers such as ViT.

## 18 1 Introduction

19 In recent years, deep learning has been extremely successful in processing massive amounts of  
20 high-dimensional and multi-modal data. Much of this success is owed to effective learning of the data  
21 distribution and then transforming the distribution to a parsimonious, i.e. *structured and compact*,  
22 representation [39, 49, 51, 61], which facilitates many downstream tasks (e.g., in vision, classification  
23 [23, 40], recognition and segmentation [25, 38, 73], and generation [31, 64, 65]). To this end, many  
24 models and methods have been proposed and practiced, each with its own strengths and limitations.  
25 Here, we give several popular methods a brief accounting as context for a complete understanding  
26 and unification that we seek in this work.

27 **Transformer models and self-attention.** Transformers [28] are one of the latest popular models  
28 for learning a representation for high-dimensional structured data, such as text [28, 30, 37], images  
29 [40, 72], and other types of signals [48, 56]. After the first block, which converts each data point  
30 (such as a text corpus or image) into a set or sequence of *tokens*, further processing is performed  
31 on the token sets, in a medium-agnostic manner [28, 40]. A cornerstone of the transformer model  
32 is the so-called *self-attention layer*, which exploits the statistical correlations among the sequence  
33 of tokens to refine the token representation. Transformers have been highly successful in learning  
34 compact representations that perform well on many downstream tasks. Yet the transformer network  
35 architecture is empirically designed and lacks a rigorous mathematical interpretation. In fact, the  
36 output of the attention layer itself has several competing interpretations [67, 74]. As a result, the  
37 statistical and geometric relationship between the data distribution and the final representation learned  
38 by a transformer largely remains a mysterious black box.

39 **Diffusion models and denoising.** Diffusion models [22, 34, 41, 43, 44] have recently become  
40 a popular method for learning the data distribution, particularly for generative tasks and natural  
41 image data which are highly structured but notoriously difficult to effectively model [3, 5]. The core  
42 concept of diffusion models is to start with features sampled from a Gaussian noise distribution (or



**Figure 1: The ‘main loop’ of the CRATE white-box deep network design.** After encoding input data  $X$  as a sequence of tokens  $Z^0$ , CRATE constructs a deep network that transforms the data to a canonical configuration of low-dimensional subspaces by successive *compression* against a local model for the distribution, generating  $Z^{\ell+1/2}$ , and *sparsification* against a global dictionary, generating  $Z^{\ell+1}$ . Repeatedly stacking these blocks and training the model parameters via backpropagation yields a powerful and interpretable representation of the data.

43 some other standard template) and *iteratively denoise* and deform the feature distribution until it  
 44 converges to the original data distribution. This process is computationally intractable if modeled in  
 45 just one step [60], so it is typically broken into multiple incremental steps. The key to each step is  
 46 the so-called *score function*, or equivalently [13] an estimate for the “optimal denoising function”;  
 47 in practice this function is modeled using a generic black-box deep network. Diffusion models  
 48 have shown effectiveness at learning and sampling from the data distribution [55, 59, 64]. However,  
 49 despite some recent efforts [77], they generally do not establish any clear correspondence between  
 50 the initial features and data samples. Hence, diffusion models themselves do not offer a parsimonious  
 51 or interpretable representation of the data distribution.

52 **Structure-seeking models and rate reduction.** In both of the previous two methods, the represen-  
 53 tations were constructed implicitly as a byproduct of solving a downstream task (e.g., classification  
 54 or generation/sampling) using deep networks. However, one can also explicitly learn a representation  
 55 of the data distribution as a task in and of itself; this is most commonly done by trying to identify and  
 56 represent low-dimensional structures in the input data. Classical examples of this paradigm include  
 57 model-based approaches such as sparse coding [2, 29] and dictionary learning [17, 21, 47], out of  
 58 which grew early attempts at designing and interpreting deep network architectures [18, 32]. More  
 59 recent approaches build instead from a model-free perspective, where one learns a representation  
 60 through a sufficiently-informative pretext task (such as compressing similar and separating dissimilar  
 61 data in contrastive learning [45, 68, 76], or maximizing the information gain in the class of maximal  
 62 coding rate reduction methods [6, 46, 54]). Compared to black-box deep learning approaches, both  
 63 model-based and model-free representation learning schemes have the advantage of being more  
 64 interpretable: they allow users to explicitly design desired properties of the learned representation [46,  
 65 54, 62]. Furthermore, they allow users to construct new white-box forward-constructed deep network  
 66 architectures [11, 54, 58] by *unrolling the optimization strategy for the representation learning*  
 67 *objective*, such that each layer of the constructed network implements an iteration of the optimization  
 68 algorithm [11, 52, 54]. Unfortunately, in this paradigm, if the desired properties are narrowly defined,  
 69 it may be difficult to achieve good practical performance on large real-world datasets.

70 **Our contributions, and outline of this work.** In this work, we aim to remedy the limitations  
 71 of these existing methods with a more unified framework for designing transformer-like network  
 72 architectures that leads to both mathematical interpretability and good practical performance. To  
 73 this end, we propose to learn a sequence of *incremental mappings* to obtain a most *compressed and*  
 74 *sparse* representation for the input data (or their token sets) that optimizes *a unified objective function*  
 75 known as the sparse rate reduction, specified later in (1). The goal of the mapping is illustrated  
 76 in Figure 1. Within this framework, we unify the above three seemingly disparate approaches and  
 77 show that *transformer-like deep network layers can be naturally derived from unrolling iterative*  
 78 *optimization schemes to incrementally optimize the sparse rate reduction objective*. In particular, our  
 79 contributions and outline of the paper are as follows:

- 80 • In Section 2.2 we show, using an idealized model for the token distribution, that if one *iteratively*  
 81 *denoises* the tokens towards a family of low-dimensional subspaces, the associated score function  
 82 assumes an explicit form similar to a self-attention operator seen in transformers.
- 83 • In Section 2.3 we derive the multi-head self-attention layer as an unrolled gradient descent step to  
 84 minimize the lossy coding rate part of the rate reduction, showing another interpretation of the  
 85 self-attention layer as compressing the token representation.
- 86 • In Section 2.4 we show that the multi-layer perceptron which immediately follows the multi-  
 87 head self-attention in transformer blocks can be interpreted as (and replaced by) a layer which

88 incrementally optimizes the remaining part of the sparse rate reduction objective by constructing  
 89 a sparse coding of the token representations.

90 • In Section 2.5 we use this understanding to create a new white-box (fully mathematically in-  
 91 terpretable) transformer architecture called CRATE (i.e., Coding RAtE reduction TransformEr),  
 92 where each layer performs a *single step* of an alternating minimization algorithm to optimize the  
 93 sparse rate reduction objective.

94 Hence, within our framework, the learning objective function, the deep learning architecture, and  
 95 the final learned representation *all become white boxes* that are fully mathematically interpretable.  
 96 As the experiments in Section 3 show, the CRATE networks, despite being simple, can already learn  
 97 the desired compressed and sparse representations on large-scale real-world datasets and achieve  
 98 performance on par with much more heavily engineered transformer networks (such as ViT) on a  
 99 wide variety of tasks (e.g., classification and transfer learning).

## 100 2 Technical Approach and Justification

### 101 2.1 Objective and Approach

102 We consider a general learning setup associated with real-world signals. We have some random  
 103 variable  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$  which is our data source; each  $\mathbf{x}_i \in \mathbb{R}^D$  is interpreted as a  
 104 *token*<sup>1</sup>, and the  $\mathbf{x}_i$ 's may have arbitrary correlation structures. We use  $\mathbf{Z} = [z_1, \dots, z_N] \in \mathbb{R}^{d \times N}$  to  
 105 denote the random variable which defines our representations. Each  $z_i \in \mathbb{R}^d$  is the representation of  
 106 the corresponding token  $\mathbf{x}_i$ . We are given  $B \geq 1$  i.i.d. samples  $\mathbf{X}_1, \dots, \mathbf{X}_B \sim \mathbf{X}$ , whose tokens are  
 107  $\mathbf{x}_{i,b}$ . The representations of our samples are denoted  $\mathbf{Z}_1, \dots, \mathbf{Z}_B \sim \mathbf{Z}$ , and those of our tokens are  
 108  $z_{i,b}$ . Finally, for a given network, we use  $\mathbf{Z}^\ell$  to denote the output of the first  $\ell$  layers when given  $\mathbf{X}$   
 109 as input. Correspondingly, the sample outputs are  $\mathbf{Z}_b^\ell$  and the token outputs are  $z_{i,b}^\ell$ .

110 **Objective for learning a structured and compact representation.** Following the framework of  
 111 rate reduction [54], we contend that the goal of representation learning is to find a feature mapping  
 112  $f: \mathbf{X} \in \mathbb{R}^{D \times N} \rightarrow \mathbf{Z} \in \mathbb{R}^{d \times N}$  which transforms input data  $\mathbf{X} \in \mathbb{R}^{D \times N}$  with a potentially  
 113 nonlinear and multi-modal distribution to a (piecewise) *linearized and compact* feature representation  
 114  $\mathbf{Z} \in \mathbb{R}^{d \times N}$ . While the joint distribution of tokens  $(z_i)_{i=1}^N$  in  $\mathbf{Z}$  may be sophisticated (and task-  
 115 specific), we further contend that it is reasonable and practical to require that the target marginal  
 116 distribution of individual tokens  $z_i$  should be highly compressed and structured, amenable for compact  
 117 coding. Particularly, we require the distribution to be a *mixture of low-dimensional (say  $K$ ) Gaussian*  
 118 *distributions*, such that the  $k^{\text{th}}$  Gaussian has mean  $\mathbf{0} \in \mathbb{R}^d$ , covariance  $\Sigma_k \succ \mathbf{0} \in \mathbb{R}^{d \times d}$ , and support  
 119 spanned by the orthonormal basis  $\mathbf{U}_k \in \mathbb{R}^{d \times p}$ . We denote  $\mathbf{U}_{[K]} = (\mathbf{U}_k)_{k=1}^K$  to be the set of bases  
 120 of all Gaussians. Hence to maximize the *information gain* [61] for the final token representation,  
 121 we wish to maximize the rate reduction [6, 46] of the tokens, i.e.,  $\max_{\mathbf{Z}} \Delta R(\mathbf{Z}; \mathbf{U}_{[K]}) = R(\mathbf{Z}) -$   
 122  $R^c(\mathbf{Z}; \mathbf{U}_{[K]})$ , where  $R$  and  $R^c$  are estimates of lossy coding rates to be formally defined in (7)  
 123 and (8). This also promotes token representations  $z_i$  from different Gaussians to be *incoherent* [46].  
 124 Since rate reduction is an intrinsic measure of goodness for the representation, it is invariant to  
 125 arbitrary rotations of the representations. Therefore, to ensure the final representations are amenable  
 126 to more compact coding, we would like to transform the representations (and their supporting  
 127 subspaces) so that they become *sparse* with respect to the standard coordinates of the resulting  
 128 representation space.<sup>2</sup> The combined rate reduction and sparsification process is illustrated in Figure 1.  
 129 Computationally, we may combine the above two goals into a unified objective for optimization:

$$\max_{f \in \mathcal{F}} \mathbb{E}_{\mathbf{Z}} [\Delta R(\mathbf{Z}; \mathbf{U}_{[K]}) - \lambda \|\mathbf{Z}\|_0] = \max_{f \in \mathcal{F}} \mathbb{E}_{\mathbf{Z}} [R(\mathbf{Z}) - R^c(\mathbf{Z}; \mathbf{U}_{[K]}) - \lambda \|\mathbf{Z}\|_0] \text{ s.t. } \mathbf{Z} = f(\mathbf{X}), \quad (1)$$

130 where the  $\ell^0$  norm  $\|\mathbf{Z}\|_0$  promotes the sparsity of the final token representations  $\mathbf{Z} = f(\mathbf{X})$ .<sup>3</sup> We  
 131 call this objective “*sparse rate reduction*.”

133 **White-box deep architecture as unrolled incremental optimization.** Although easy to state, each  
 134 term of the above objective can be computationally very challenging to optimize [54, 69]. Hence it is  
 135 natural to take an approximation approach that realizes the global transformation  $f$  optimizing (1)  
 136 through a concatenation of multiple, say  $L$ , simple *incremental and local* operations  $f^\ell$  that push the  
 137 representation distribution towards the desired parsimonious model distribution:

<sup>1</sup>For language transformers, tokens roughly correspond to words [28], while for vision transformers, tokens correspond to image patches [40].

<sup>2</sup>That is, having the fewest nonzero entries.

<sup>3</sup>To simplify the notation, we will discuss the objective for one sample  $\mathbf{X}$  at a time with the understanding that we always mean to optimize the expectation.

$$f: \mathbf{X} \xrightarrow{f^0} \mathbf{Z}^0 \rightarrow \dots \rightarrow \mathbf{Z}^\ell \xrightarrow{f^\ell} \mathbf{Z}^{\ell+1} \rightarrow \dots \rightarrow \mathbf{Z}^L = \mathbf{Z}, \quad (2)$$

138 where  $f^0: \mathbb{R}^D \rightarrow \mathbb{R}^d$  is the pre-processing mapping that transforms input tokens  $\mathbf{x}_i \in \mathbb{R}^D$  to their  
 139 token representations  $\mathbf{z}_i^1 \in \mathbb{R}^d$ . Each incremental *forward mapping*  $\mathbf{Z}^{\ell+1} = f^\ell(\mathbf{Z}^\ell)$ , or a “layer”,  
 140 transforms the token distribution to *optimize* the above sparse rate reduction objective (1), conditioned  
 141 on the distribution of its input tokens  $\mathbf{Z}^\ell$ . The distribution of  $\mathbf{Z}^\ell$  can be explicitly modeled or  
 142 approximated, say as a mixture of linear subspaces or sparsely generated from a dictionary, with  
 143 parameters learned from data (say via *backward propagation* with end-to-end training).<sup>4</sup>

144 We show that we can derive these incremental, local operations through an unrolled optimization  
 145 perspective to achieve (1) through Sections 2.3 to 2.5. Once we decide on using an incremental  
 146 approach to optimizing (1), there are a variety of possible choices to achieve the optimization. Given  
 147 a model for  $\mathbf{Z}^\ell$ , say a mixture of subspaces  $\mathbf{U}_{[K]}$ , we opt for a two-step *alternating minimization*  
 148 process with a strong conceptual basis: first in Section 2.3, we *compress* the tokens  $\mathbf{Z}^\ell$  via a gradient  
 149 step to minimize the coding rate term  $\min_{\mathbf{Z}} R^c(\mathbf{Z}; \mathbf{U}_{[K]})$ ; second, in Section 2.4, we *sparsify* the  
 150 compressed tokens, with a suitably-relaxed proximal gradient step on the difference of the sparsity  
 151 penalty and the expansion term, i.e.,  $\min_{\mathbf{Z}} [\lambda \|\mathbf{Z}\|_0 - R(\mathbf{Z})]$ . Both actions are applied incrementally  
 152 and repeatedly, as each  $f^\ell$  in (2) is instantiated with these two steps.

## 153 2.2 Self-Attention via Denoising Tokens Towards Multiple Subspaces

154 There are many different ways to optimize the objective (1) incrementally. In this work, we propose  
 155 arguably *the most basic* scheme. To help clarify the intuition behind our derivation and approximation,  
 156 in this section (and Appendix A.1) we study a largely idealized model which nevertheless captures  
 157 the essence of nearly the whole process and particularly reveals the reason why self-attention-like  
 158 operators arise in many contexts. Assume that  $N = 1$ , and the single token  $\mathbf{x}$  is drawn i.i.d. from  
 159 an unknown mixture of Gaussians  $(\mathcal{N}(\mathbf{0}, \Sigma_k))_{k=1}^K$  supported on low-dimensional subspaces with  
 160 orthonormal bases  $\mathbf{U}_{[K]} = (\mathbf{U}_k)_{k=1}^K$  and corrupted with additive Gaussian noise  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , i.e.,

$$\mathbf{x} = \mathbf{z} + \sigma \mathbf{w}, \quad (3)$$

161 where  $\mathbf{z}$  is distributed according to the mixture. Our goal is simply to transform the distribution of  
 162 the noisy token  $\mathbf{x}$  to the mixture of low-dimensional Gaussians  $\mathbf{z}$ . Towards incremental construction  
 163 of a representation  $f$  for this model following (2), we reason inductively: if  $\mathbf{z}^\ell$  is a noisy token (3)  
 164 at noise level  $\sigma^\ell$ , it is natural to produce  $\mathbf{z}^{\ell+1}$  by denoising at the level  $\sigma^\ell$ . In the mean-square sense,  
 165 the optimal estimate is  $\mathbb{E}[\mathbf{z} | \mathbf{z}^\ell]$ , which has a variational characterization (e.g. [12]):

$$\mathbb{E}[\mathbf{z} | \cdot] = \arg \min_f \mathbb{E}_{\mathbf{z}, \mathbf{w}} \left[ \left\| f(\mathbf{z} + \sigma^\ell \mathbf{w}) - \mathbf{z} \right\|_2^2 \right]. \quad (4)$$

166 Setting  $\mathbf{z}^{\ell+1} = \mathbb{E}[\mathbf{z} | \mathbf{z}^\ell]$ , (4) thus characterizes the next stage of (2) in terms of an optimization  
 167 objective based on a *local signal model* for  $\mathbf{z}^\ell$ . Moreover, letting  $\mathbf{x} \mapsto q^\ell(\mathbf{x})$  denote the density of  $\mathbf{z}^\ell$ ,  
 168 Tweedie’s formula [13] allows us to express the optimal representation solving (4) in closed-form:

$$\mathbf{z}^{\ell+1} = \mathbf{z}^\ell + (\sigma^\ell)^2 \nabla_{\mathbf{x}} \log q^\ell(\mathbf{z}^\ell). \quad (5)$$

169 Tweedie’s formula expresses the optimal representation in terms of an additive correction (in general  
 170 a nonlinear function of  $\mathbf{z}^\ell$ ) to the noisy observations by the gradient of the *log-likelihood* of the  
 171 distribution of the noisy observations, giving the optimal representation a clear interpretation as an  
 172 incremental perturbation to the current noisy distribution  $q^\ell$ . This connection is well-known in the  
 173 areas of estimation theory and inverse problems [1, 13, 14, 19, 20, 27, 42], and more recently has  
 174 found powerful applications in the training of generative models for natural images [4, 15, 22, 43,  
 175 44]. Here, we can calculate a closed-form expression for this *score function*  $\nabla_{\mathbf{x}} \log q^\ell$ , which, when  
 176 combined with (5) and some technical assumptions<sup>5</sup>, gives the following approximation (shown in  
 177 Appendix A.1). Let  $\otimes$  denote the Kronecker product; then we have

$$\mathbf{z}^{\ell+1} \approx [\mathbf{U}_1, \dots, \mathbf{U}_K] \left[ \text{diag} \left( \text{softmax} \left( \frac{1}{2(\sigma^\ell)^2} \begin{bmatrix} \|\mathbf{U}_1^* \mathbf{z}^\ell\|_2^2 \\ \vdots \\ \|\mathbf{U}_K^* \mathbf{z}^\ell\|_2^2 \end{bmatrix} \right) \right) \otimes \mathbf{I}_p \right] \begin{bmatrix} \mathbf{U}_1^* \mathbf{z}^\ell \\ \vdots \\ \mathbf{U}_K^* \mathbf{z}^\ell \end{bmatrix}, \quad (6)$$

<sup>4</sup>This separation of forward “optimization” and backward “learning” clarifies the mathematical role of each layer as an operator transforming the distribution of its input, whereas the input distribution is in turn modeled by the parameters of the layer.

<sup>5</sup>Such as  $\sigma$  being smaller than the nonzero eigenvalues of  $\Sigma_k$  and the normalization assumption  $\pi_i \det(\Sigma_i + \sigma^2 \mathbf{I})^{-1/2} = \pi_j \det(\Sigma_j + \sigma^2 \mathbf{I})^{-1/2}$  for all  $i, j \in [K]$ , where  $\pi_k$  is the mixture proportion for the  $k^{\text{th}}$  Gaussian.

178 This operation resembles a self-attention layer in a standard transformer architecture with  $K$  heads,  
 179 sequence length  $N = 1$ , the “query-key-value” constructs being replaced by a single linear projection  
 180  $\mathbf{U}_k^* \mathbf{z}^\ell$  of the token  $\mathbf{z}^\ell$ , and the aggregation of head outputs (conventionally modeled by an MLP)  
 181 done with the two leftmost matrices in (6). We thus derive the following useful interpretation, which  
 182 we will exploit in the sequel: *Gaussian denoising against a mixture of subspaces model leads to*  
 183 *self-attention-type layers in the transformation  $f$ .* Given an initial sample  $\mathbf{x}$  following the model  
 184 (3), we can repeatedly apply local transformations to the distribution with (6) in order to realize the  
 185 incremental mapping  $f: \mathbf{x} \rightarrow \mathbf{z}$  in (2).<sup>6</sup> These insights will guide us in the design of our white-box  
 186 transformer architecture in the upcoming subsections.

### 187 2.3 Self-Attention via Compressing Token Sets through Optimizing Rate Reduction

188 In the last subsection, we have seen that the multi-head attention in a transformer resembles the score-  
 189 matching operator that aims to transform a token  $\mathbf{z}^\ell$  towards a mixture of subspaces (or degenerate  
 190 Gaussians). Nevertheless, to carry out such an operation on any data, one needs to first learn or  
 191 estimate, typically from finite samples, the parameters of the mixture of (degenerate) Gaussians,  
 192 which is known to be a challenging task [6, 24]. This challenge is made even harder because in a  
 193 typical learning setting, the given set of tokens are *not* i.i.d. samples from the mixture of subspaces.  
 194 The joint distribution among these tokens can encode rich information about the data—for example,  
 195 co-occurrences between words or object parts in language and image data (resp.)—which we should  
 196 also learn. Thus, we should compress/denoise/transform such a set of tokens together. To this end,  
 197 we need a measure of quality, i.e., compactness, for the resulting representation of the set of tokens.  
 198 A natural measure of the compactness of such a set of tokens is the (lossy) coding rate to encode  
 199 them up to a certain precision  $\epsilon > 0$  [6, 46]. For a zero-mean Gaussian, this measure takes a closed  
 200 form. If we view the tokens in  $\mathbf{Z} \in \mathbb{R}^{d \times N}$  as drawn from a single zero-mean Gaussian, an estimate  
 201 of their (lossy) coding rate, subject to quantization precision  $\epsilon > 0$ , is given in [6] as:

$$R(\mathbf{Z}) \doteq \frac{1}{2} \log \det \left( \mathbf{I} + \frac{d}{N\epsilon^2} \mathbf{Z}^* \mathbf{Z} \right) = \frac{1}{2} \log \det \left( \mathbf{I} + \frac{d}{N\epsilon^2} \mathbf{Z} \mathbf{Z}^* \right). \quad (7)$$

202 In practice, the data distribution is typically multi-modal, say an image set consisting of many classes  
 203 or a collection of image patches as in Figure 1. It is more appropriate to require that the set of  
 204 tokens map to a mixture of, say  $K$ , subspaces (degenerate Gaussians) [54]. As before we denote  
 205 the (to be learned) bases of these subspaces as  $\mathbf{U}_{[K]} = (\mathbf{U}_k)_{k=1}^K$ , where  $\mathbf{U}_k \in \mathbb{R}^{d \times p}$ . Although the  
 206 joint distribution of the tokens  $\mathbf{Z}$  is unknown, the desired marginal distribution of each token  $\mathbf{z}_i$  is a  
 207 mixture of subspaces. So we may obtain an upper bound of the coding rate for the token set  $\mathbf{Z}$  by  
 208 projecting its tokens onto these subspaces and summing up the respective coding rates:

$$R^c(\mathbf{Z}; \mathbf{U}_{[K]}) = \sum_{k=1}^K R(\mathbf{U}_k^* \mathbf{Z}) = \frac{1}{2} \sum_{k=1}^K \log \det \left( \mathbf{I} + \frac{p}{N\epsilon^2} (\mathbf{U}_k^* \mathbf{Z})^* (\mathbf{U}_k^* \mathbf{Z}) \right). \quad (8)$$

209 We would like to compress (or denoise) the set of tokens against these subspaces by minimizing the  
 210 coding rate. The gradient of  $R^c(\mathbf{Z}; \mathbf{U}_{[K]})$  is

$$\nabla_{\mathbf{Z}} R^c(\mathbf{Z}; \mathbf{U}_{[K]}) = \frac{p}{N\epsilon^2} \sum_{k=1}^K \mathbf{U}_k \mathbf{U}_k^* \mathbf{Z} \left( \mathbf{I} + \frac{p}{N\epsilon^2} (\mathbf{U}_k^* \mathbf{Z})^* (\mathbf{U}_k^* \mathbf{Z}) \right)^{-1}. \quad (9)$$

211 The above expression approximates the residual of each projected token  $\mathbf{U}_k^* \mathbf{z}_i$  regressed by other  
 212 tokens  $\mathbf{U}_k^* \mathbf{z}_j$  [54]. But, differently from [54], not all tokens in  $\mathbf{Z}$  are from the same subspace. Hence,  
 213 to denoise each token with tokens from its own group, we can compute their similarity through an  
 214 auto-correlation among the projected tokens as  $(\mathbf{U}_k^* \mathbf{Z})^* (\mathbf{U}_k^* \mathbf{Z})$  and convert it to a distribution of  
 215 membership with a softmax, namely  $\text{softmax}((\mathbf{U}_k^* \mathbf{Z})^* (\mathbf{U}_k^* \mathbf{Z}))$ . Then, as we show in Appendix A.2,  
 216 if we only use similar tokens to regress and denoise each other, then a gradient step on the coding  
 217 rate with learning rate  $\kappa$  can be naturally approximated as follows:

$$\mathbf{Z}^{\ell+1/2} = \mathbf{Z}^\ell - \kappa \nabla_{\mathbf{Z}} R^c(\mathbf{Z}^\ell; \mathbf{U}_{[K]}) \approx \left( 1 - \kappa \cdot \frac{p}{N\epsilon^2} \right) \mathbf{Z}^\ell + \kappa \cdot \frac{p}{N\epsilon^2} \cdot \text{MSSA}(\mathbf{Z}^\ell \mid \mathbf{U}_{[K]}), \quad (10)$$

218 where MSSA is defined through an SSA operator as:

$$\text{SSA}(\mathbf{Z} \mid \mathbf{U}_k) \doteq (\mathbf{U}_k^* \mathbf{Z}) \text{softmax}((\mathbf{U}_k^* \mathbf{Z})^* (\mathbf{U}_k^* \mathbf{Z})), \quad k \in [K], \quad (11)$$

<sup>6</sup>This statement can be made mathematically rigorous by exploiting a deep connection between neural ODEs and diffusion models, following ideas in Song et al. [44] and Chen et al. [70].

$$\text{MSSA}(\mathbf{Z} \mid \mathbf{U}_{[K]}) \doteq \frac{p}{N\epsilon^2} \cdot [\mathbf{U}_1, \dots, \mathbf{U}_K] \begin{bmatrix} \text{SSA}(\mathbf{Z} \mid \mathbf{U}_1) \\ \vdots \\ \text{SSA}(\mathbf{Z} \mid \mathbf{U}_K) \end{bmatrix}. \quad (12)$$

219 Here the SSA operator in (11) resembles the *attention operator* in a typical transformer [28], except  
 220 that here the linear operators of value, key, and query are all set to be *the same* as the subspace  
 221 basis, i.e.,  $\mathbf{V} = \mathbf{K} = \mathbf{Q} = \mathbf{U}_k^*$ .<sup>7</sup> Hence, we name  $\text{SSA}(\cdot \mid \mathbf{U}_k) : \mathbb{R}^{d \times N} \rightarrow \mathbb{R}^{p \times N}$  the **Subspace**  
 222 **Self-Attention (SSA)** operator (more details and justification can be found in (71) in Appendix A.2).  
 223 Then, the whole MSSA operator in (12), formally defined as  $\text{MSSA}(\cdot \mid \mathbf{U}_{[K]}) : \mathbb{R}^{d \times N} \rightarrow \mathbb{R}^{d \times N}$  and  
 224 called the **Multi-Head Subspace Self-Attention (MSSA)** operator, aggregates the attention head  
 225 outputs by averaging using model-dependent weights, similar in concept to the popular multi-head  
 226 self-attention operator in existing transformer networks. The overall gradient step (10) resembles the  
 227 multi-head self-attention implemented with a skip connection in transformers.

228 Notice that if we have  $N = 1$  tokens as well as take an aggressive gradient step ( $\kappa = 1$ ) and tune the  
 229 quantization error ( $\epsilon = \sqrt{p/N}$ ), the multi-head subspace self-attention operator in (12) becomes the  
 230 ideal denoiser defined in (6), with the one minor difference that the aggregation of the heads is done  
 231 by a linear function here, while in (6) it is done by a nonlinear mixture-of-experts type function.<sup>8</sup>  
 232 This provides two very related interpretations of the multi-head self-attention operator, as denoising  
 233 and compression against a mixture of low-dimensional subspaces.

## 2.4 MLP via Iterative Shrinkage-Thresholding Algorithms (ISTA) for Sparse Coding

235 In the previous subsection, we focused on how to compress a set of tokens against a set of (learned)  
 236 low-dimensional subspaces. Optimizing the remaining terms in the sparse rate reduction objective  
 237 (1), including the non-smooth term, serves to sparsify the compressed tokens, hence leading to a  
 238 more compact and structured (i.e., *parsimonious*) representation. From (1) and (7), this term is

$$\max_{\mathbf{Z}} [R(\mathbf{Z}) - \lambda \|\mathbf{Z}\|_0] = \min_{\mathbf{Z}} \left[ \lambda \|\mathbf{Z}\|_0 - \frac{1}{2} \log \det \left( \mathbf{I} + \frac{d}{N\epsilon^2} \mathbf{Z}^* \mathbf{Z} \right) \right], \quad (13)$$

239 where  $R(\mathbf{Z})$  denotes the coding rate of the whole token set, as defined in (7). In addition to  
 240 sparsification via the  $\|\mathbf{Z}\|_0$  term, the expansion term  $R(\mathbf{Z})$  in (13) promotes diversity and non-  
 241 collapse of the representation, a highly desirable property. However, prior work has struggled to  
 242 realize this benefit on large-scale datasets due to poor scalability of the gradient  $\nabla_{\mathbf{Z}} R(\mathbf{Z})$ , which  
 243 requires a matrix inverse [54].

244 To simplify things, we therefore take a different approach to trading off between representational  
 245 diversity and sparsification: we posit a (complete) incoherent or orthogonal dictionary  $\mathbf{D} \in \mathbb{R}^{d \times d}$ , and  
 246 ask to sparsify the intermediate iterates  $\mathbf{Z}^{\ell+1/2}$  with respect to  $\mathbf{D}$ . That is,  $\mathbf{Z}^{\ell+1/2} = \mathbf{D}\mathbf{Z}^{\ell+1}$  where  
 247  $\mathbf{Z}^{\ell+1}$  is more sparse. The dictionary  $\mathbf{D}$  is global, i.e., is used to sparsify all tokens simultaneously.  
 248 By the incoherence assumption, we have  $\mathbf{D}^* \mathbf{D} \approx \mathbf{I}_d$ ; thus from (7) we have  $R(\mathbf{Z}^{\ell+1}) \approx$   
 249  $R(\mathbf{D}\mathbf{Z}^{\ell+1}) = R(\mathbf{Z}^{\ell+1/2})$ . Thus we approximately solve (13) with the following program:

$$\min_{\mathbf{Z}^{\ell+1}} \|\mathbf{Z}^{\ell+1}\|_0 \quad \text{subject to} \quad \mathbf{Z}^{\ell+1/2} = \mathbf{D}\mathbf{Z}^{\ell+1}. \quad (14)$$

250 The above sparse representation program is usually solved by relaxing it to an unconstrained convex  
 251 program, known as LASSO:  $\min_{\mathbf{Z}^{\ell+1}} [\lambda \|\mathbf{Z}^{\ell+1}\|_1 + \|\mathbf{Z}^{\ell+1/2} - \mathbf{D}\mathbf{Z}^{\ell+1}\|_F^2]$ . In our implementation,  
 252 motivated by Sun et al. [33] and Zarka et al. [35], we also add a non-negative constraint to  $\mathbf{Z}^{\ell+1}$ ,

$$\mathbf{Z}^{\ell+1} = \arg \min_{\mathbf{Z} \geq 0} [\lambda \|\mathbf{Z}\|_1 + \|\mathbf{Z}^{\ell+1/2} - \mathbf{D}\mathbf{Z}\|_F^2], \quad (15)$$

253 which we then incrementally optimize by performing an unrolled proximal gradient descent step,  
 254 known as an ISTA step [8], to give the update:

$$\mathbf{Z}^{\ell+1} = \text{ReLU}(\mathbf{Z}^{\ell+1/2} + \eta \mathbf{D}^* (\mathbf{Z}^{\ell+1/2} - \mathbf{D}\mathbf{Z}^{\ell+1/2}) - \eta \lambda \mathbf{1}) \doteq \text{ISTA}(\mathbf{Z}^{\ell+1/2} \mid \mathbf{D}). \quad (16)$$

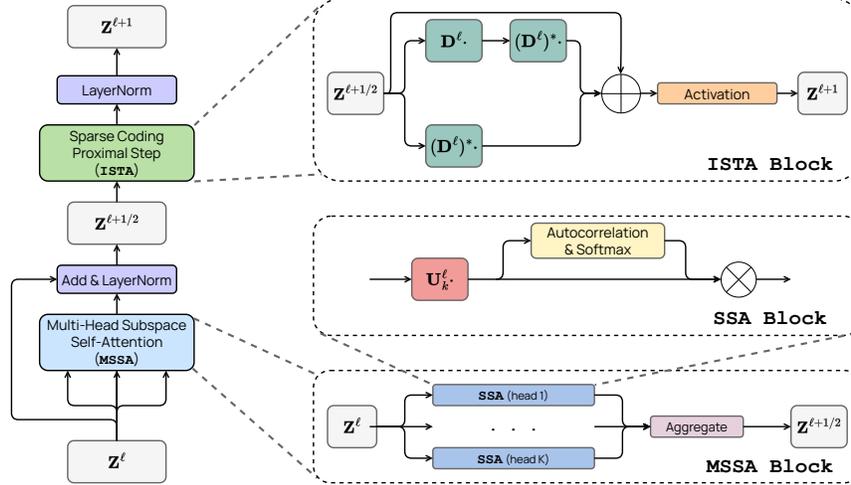
255 In Appendix A.3, we will show one can arrive at a similar operator to the above ISTA-like update for  
 256 optimizing (13) by properly linearizing and approximating the rate term  $R(\mathbf{Z})$ .

## 2.5 The Overall White-Box CRATE Architecture

258 By combining the above two steps:

<sup>7</sup>We note a recent suggestion of Hinton [50] that it is more sensible to set the “value, key, and query” projection matrices in a transformer to be equal. Our derivation in this section confirms this mathematically.

<sup>8</sup>This suggests that we could also consider such a mixture of expert type aggregation of the multiple attention heads. In this work, we use linear aggregation, and leave evaluation of more variants for future work.



**Figure 2:** One layer of the CRATE architecture. The full architecture is simply a concatenation of such layers, with some initial tokenizer and final task-specific architecture (i.e., a classification head).

- 259 1. (Sections 2.2 and 2.3) Local denoising and compression of tokens within a sample towards a  
 260 mixture-of-subspace structure, leading to the multi-head subspace self-attention block – MSSA;  
 261 2. (Section 2.4) Global compression and sparsification of token sets across all samples through  
 262 sparse coding, leading to the sparsification block – ISTA;  
 263 we can get the following rate-reduction-based transformer layer, illustrated in Figure 2,

$$\mathbf{Z}^{\ell+1/2} \doteq \mathbf{Z}^{\ell} + \text{MSSA}(\mathbf{Z}^{\ell} \mid \mathbf{U}_{[K]}^{\ell}), \quad \mathbf{Z}^{\ell+1} \doteq \text{ISTA}(\mathbf{Z}^{\ell+1/2} \mid \mathbf{D}^{\ell}). \quad (17)$$

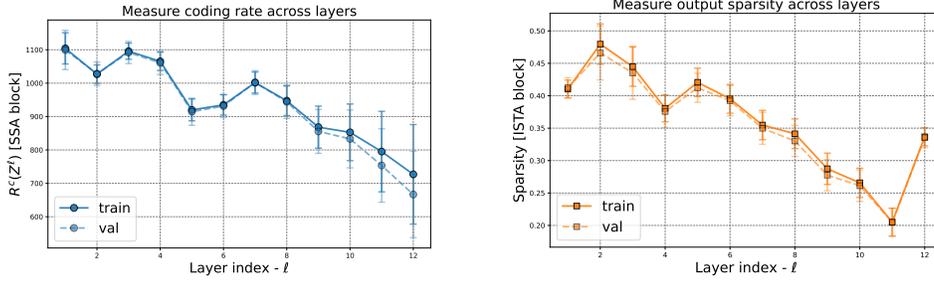
264 Composing multiple such layers following the incremental construction of our representation in (2),  
 265 we obtain a white-box transformer architecture that transforms the data tokens towards a compact  
 266 and sparse union of incoherent subspaces.

267 This model has the parameters  $(\mathbf{U}_{[K]}^{\ell})_{\ell=1}^L$  and  $(\mathbf{D}^{\ell})_{\ell=1}^L$ , which are learned from data via *back-*  
 268 *propagation*. Notably, in each layer  $\ell$ , the learned  $\mathbf{U}_{[K]}^{\ell}$  retain their interpretation as incoherent bases  
 269 for supporting subspaces for the mixture-of-Gaussians model at layer  $\ell$ , and the learned  $\mathbf{D}^{\ell}$  retains its  
 270 interpretation as a sparsifying dictionary at layer  $\ell$ . The parameters depend on the layer  $\ell$  so as to adapt  
 271 to local properties of the data distribution at each layer of the network. Our interpretation clarifies the  
 272 roles of the network forward pass (given local signal models at each layer, denoise/compress/sparsify  
 273 the input) and the backward pass (learn the local signal models from data).

274 We note that at each stage of our construction, we have chosen the *simplest possible* construction  
 275 to use. We can substitute each part of this construction, so long as the new part maintains the same  
 276 conceptual role, and obtain another white-box architecture. Nevertheless, our such-constructed  
 277 architecture, called CRATE (i.e., Coding RAtE TransformEr), connects to existing transformer models,  
 278 obtains competitive results on real-world datasets, and is fully mathematically interpretable.

### 279 3 Experiments

280 In this section, we conduct experiments to study the performance of our proposed white-box trans-  
 281 former CRATE on real-world datasets and tasks. As the analysis in Section 2 suggests, either the  
 282 compression or the sparsification step can be achieved through various alternative design choices or  
 283 strategies. CRATE arguably adopts the most basic choices and so our goal with the experiments is *not*  
 284 simply to compete with other heavily engineered transformers while using such a rudimentary design.  
 285 Rather, our goals are twofold. First, unlike any empirically designed black-box networks that are  
 286 usually evaluated only on end-to-end performance, the white-box design of our network allows us  
 287 to *look inside* the deep architecture and verify if layers of the learned network indeed perform their  
 288 design objective—say performing incremental optimization for the objective (1). Second, despite their  
 289 simplicity, our experiments will actually reveal the vast practical potential of our so-derived CRATE  
 290 architectures since, as we will show, they already achieve very strong performance on large-scale  
 291 real-world datasets and tasks. In the remainder of this section we highlight a selection of results;  
 292 additional experimental details and results can be found in Appendix B.



**Figure 3:** *Left:* The compression term  $R^c(\mathbf{Z}^{\ell+1/2})$  of the MSSA outputs at different layers. *Right:* the sparsity of the ISTA output block,  $\|\mathbf{Z}^{\ell+1}\|_0/(d \cdot N)$ , at different layers. (Model: CRATE-Small).

293 **Model architecture.** We implement the architecture that is described in Section 2.5, with minor  
 294 modifications that are described in Appendix B.1. We consider different model sizes of CRATE by  
 295 varying the token dimension  $d$ , number of heads  $K$ , and the number of layers  $L$ . We consider four  
 296 model sizes in this work: CRATE-Tiny, CRATE-Small, CRATE-Base, and CRATE-Large. A PyTorch-  
 297 style pseudocode can be found in Appendix B.1, which contains more implementation details. For  
 298 training using supervised classification, we first take the CLS token  $\bar{z}_b = z_{1,b}^{L+1}$  of for each sample,  
 299 then apply a linear layer; the output of this linear layer  $u_b \doteq \mathbf{W}\bar{z}_b$  is used as input to the standard  
 300 cross-entropy loss. The overall loss averages over all samples  $b \in [B]$ .

301 **Datasets and optimization.** We mainly consider ImageNet-1K [9] as the testbed for our architecture.  
 302 Specifically, we apply the Lion optimizer [71] to train CRATE models with different model sizes.  
 303 Meanwhile, we also evaluate the transfer learning performance of CRATE: by considering the models  
 304 trained on ImageNet-1K as pre-trained models, we fine-tune CRATE on several commonly used  
 305 downstream datasets (CIFAR10/100, Oxford Flowers, Oxford-IIT-Pets). More details about the  
 306 training and datasets can be found in Appendix B.1.

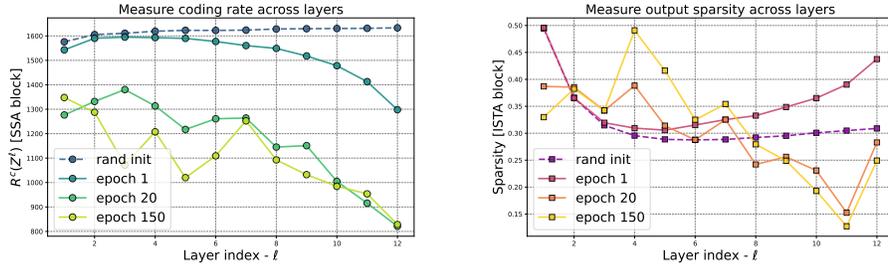
### 307 3.1 In-depth Layer-wise Analysis of CRATE

308 **Do layers of CRATE achieve their design goals?** As described in Section 2.3 and Section 2.4, the  
 309 MSSA block is designed to optimize the compression term  $R^c(\mathbf{Z})$  and the ISTA block to sparsify the  
 310 token representations (corresponding to the sparsification term  $\|\mathbf{Z}\|_0$ ). To understand whether CRATE  
 311 indeed optimizes these terms, for each layer  $\ell$ , we measure (i) the compression term  $R^c(\mathbf{Z}^{\ell+1/2})$   
 312 on the MSSA block outputs  $\mathbf{Z}^{\ell+1/2}$ ; and (ii) sparsity  $\|\mathbf{Z}^{\ell+1}\|_0$  on the ISTA block outputs  $\mathbf{Z}^{\ell+1}$ .  
 313 Specifically, we evaluate these two terms by using training/validation samples from ImageNet-1K.  
 314 Both terms are evaluated at the per-sample level and averaged over  $B = 10^3$  samples.

315 Figure 3 shows the plots of these two key measures at all layers for the learned CRATE-small model.  
 316 We find that as the layer index  $\ell$  increases, both the compression and the sparsification terms improve  
 317 in most cases. The increase in the sparsity measure of the last layer is caused by the extra linear  
 318 layer for classification. These results suggest that CRATE aligns well with the original design goals:  
 319 once learned, it essentially learns to gradually compress and sparsity the representations through  
 320 its layers. In addition, we also measure the compression and sparsification terms on CRATE models  
 321 with different model sizes as well as intermediate model checkpoints and the results are shown by  
 322 plots in Figure 5 of Appendix B.2. The observations are very consistent across all different model  
 323 sizes—both the compression and sparsification terms improve in most scenarios. Models with more  
 324 layers tend to optimize the objectives more effectively, confirming our understanding of each layer’s  
 325 roles.

326 To see the effect of learning, we present the evaluations on CRATE-Small trained with different number  
 327 of epochs in Figure 4. When the model is not trained enough (e.g. untrained), the architecture does  
 328 not optimize the objectives effectively. However, during training—learning better subspaces  $\mathbf{U}_{[K]}^\ell$   
 329 and dictionaries  $\mathbf{D}^\ell$ —the designed blocks start to optimize the objectives much more effectively.

330 **Visualizing layer-wise token representations.** To gain a better understanding of the token represen-  
 331 tations of CRATE, we visualize the output of each ISTA block at layer  $\ell$  in Figure 6 of Appendix B.2.  
 332 Specifically, we visualize the  $\mathbf{Z}^{\ell+1}$  via heatmap plots. We observe that the output  $\mathbf{Z}^{\ell+1}$  becomes  
 333 more sparse as the layer increases. Moreover, besides the sparsity, we also find that  $\mathbf{Z}^{\ell+1}$  becomes  
 334 more structured (i.e., low-rank), which indicates that the set of token representations become closer  
 335 to linear subspaces, confirming our mental picture of the geometry of each layer (as in Figure 1).



**Figure 4:** The compression term  $R^c(\mathbf{Z})$  (left) and sparsification term  $\|\mathbf{Z}\|_0/(d \cdot N)$  (right) across models trained with different numbers of epochs. (Model: CRATE-Base).

**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

336 **Visualizing layer-wise subspaces in multi-head self-attention.** We now visualize the  $U_{[K]}^\ell$  matrices  
 337 used in the MSSA block. In Section 2.3, we assumed that  $U_{[K]}^\ell$  were incoherent to capture different  
 338 “views” of the set of tokens. In Fig. 7 of Appendix B.2, we first normalize the columns in each  
 339  $U_k^\ell$ , then we visualize the  $[U_1^\ell, \dots, U_K^\ell]^* [U_1^\ell, \dots, U_K^\ell] \in \mathbb{R}^{pK \times pK}$ . The  $(i, j)$ -th block in each sub-  
 340 figure corresponds to  $(U_i^\ell)^* U_j^\ell$  for  $i, j \in [K]$  at a particular layer  $\ell$ . We find that the learned  $U_{[K]}^\ell$   
 341 are approximately incoherent, which aligns well with our assumptions. One interesting observation is  
 342 that the  $U_{[K]}^\ell$  becomes more incoherent when the layer index  $\ell$  is larger, which suggests that the token  
 343 representations are more separable. This mirrors the situation in other popular deep networks [57].

### 344 3.2 Evaluations of CRATE on Large Real-World Datasets and Tasks

345 We now study the empirical performance of the proposed networks by measuring their top-1 accuracy  
 346 on ImageNet-1K as well as transfer learning performance on several widely used downstream datasets.  
 347 We summarize the results in Table 1. As our designed architecture leverages parameter sharing in  
 348 both the attention block (MSSA) and the MLP block (ISTA), our CRATE-Base model (22.08 million)  
 349 has a similar number of parameters to the ViT-Small (22.05 million).

350 From Table 1, we find that with a similar number of model parameters, our proposed network  
 351 achieves similar ImageNet-1K and transfer learning performance as ViT, despite the simplicity and  
 352 interpretability of our design. Moreover, with the same set of training hyperparameters, we observe  
 353 promising scaling behavior in CRATE—we consistently improve the performance by scaling up the  
 354 model size. For comparison, directly scaling ViT on ImageNet-1K does not always lead to consistent  
 355 performance improvement measured by top-1 accuracy [40]. To summarize, we achieve promising  
 356 performance on real-world large-scale datasets by directly implementing our principled architecture.

## 357 4 Conclusion

358 In this paper, we propose a new theoretical framework that allows us to derive deep transformer-  
 359 like network architectures as incremental optimization schemes to learn compressed and sparse  
 360 representation of the input data (or token sets). The so derived and learned deep architectures are not  
 361 only fully mathematically interpretable, but also consistent on a layer-by-layer level with their design  
 362 objective. Despite being arguably the simplest among all possible designs, these networks already  
 363 demonstrate performance on large-scale real-world datasets and tasks close to seasoned transformers.  
 364 We believe this work truly helps bridge the gap between theory and practice of deep neural networks  
 365 as well as help unify seemingly separate approaches to learning and representing data distributions.  
 366 Probably more importantly for practitioners, our framework provides theoretical guidelines to design  
 367 and justify new, potentially more powerful, deep architectures for representation learning.

368 **References**

- 369 [1] Charles M Stein. “Estimation of the Mean of a Multivariate Normal Distribution”. *The Annals*  
370 *of Statistics* 9.6 (Nov. 1981), pp. 1135–1151. 4.
- 371 [2] Bruno A Olshausen and David J Field. “Sparse coding with an overcomplete basis set: A  
372 strategy employed by V1?” *Vision research* 37.23 (1997), pp. 3311–3325. 2.
- 373 [3] David L Donoho and Carrie Grimes. “Image Manifolds which are Isometric to Euclidean  
374 Space”. *Journal of mathematical imaging and vision* 23.1 (July 2005), pp. 5–24. 1.
- 375 [4] Aapo Hyvärinen. “Estimation of Non-Normalized Statistical Models by Score Matching”.  
376 *Journal of machine learning research: JMLR* 6.24 (2005), pp. 695–709. 4.
- 377 [5] Michael B Wakin, David L Donoho, Hyeokho Choi, and Richard G Baraniuk. “The multiscale  
378 structure of non-differentiable image manifolds”. *Wavelets XI*. Vol. 5914. SPIE. 2005, pp. 413–  
379 429. 1.
- 380 [6] Yi Ma, Harm Derksen, Wei Hong, and John Wright. “Segmentation of multivariate mixed data  
381 via lossy data coding and compression”. *PAMI* (2007). 2, 3, 5.
- 382 [7] Maria-Elena Nilsback and Andrew Zisserman. “Automated flower classification over a large  
383 number of classes”. *2008 Sixth Indian Conference on Computer Vision, Graphics & Image*  
384 *Processing*. IEEE. 2008, pp. 722–729. 24.
- 385 [8] Amir Beck and Marc Teboulle. “A fast iterative shrinkage-thresholding algorithm for linear  
386 inverse problems”. *SIAM journal on imaging sciences* 2.1 (2009), pp. 183–202. 6.
- 387 [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. “Imagenet: A large-  
388 scale hierarchical image database”. *2009 IEEE conference on computer vision and pattern*  
389 *recognition*. Ieee. 2009, pp. 248–255. 8.
- 390 [10] Alex Krizhevsky, Geoffrey Hinton, et al. “Learning multiple layers of features from tiny  
391 images” (2009). 24.
- 392 [11] Karol Gregor and Yann LeCun. “Learning fast approximations of sparse coding”. *Proceedings*  
393 *of the 27th International Conference on International Conference on Machine Learning*.  
394 Omnipress. 2010, pp. 399–406. 2.
- 395 [12] László Györfi, Michael Kohler, Adam Krzyzak, and Harro Walk. *A Distribution-Free Theory*  
396 *of Nonparametric Regression*. Springer New York, Dec. 2010. 4.
- 397 [13] Bradley Efron. “Tweedie’s Formula and Selection Bias”. *Journal of the American Statistical*  
398 *Association* 106.496 (2011), pp. 1602–1614. 2, 4, 14.
- 399 [14] Martin Raphan and Eero P Simoncelli. “Least squares estimation without priors or supervision”.  
400 *Neural computation* 23.2 (Feb. 2011), pp. 374–420. 4.
- 401 [15] Pascal Vincent. “A connection between score matching and denoising autoencoders”. *Neural*  
402 *computation* 23.7 (July 2011), pp. 1661–1674. 4.
- 403 [16] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. “Cats and dogs”. *2012*  
404 *IEEE conference on computer vision and pattern recognition*. IEEE. 2012, pp. 3498–3505. 24.
- 405 [17] Daniel A Spielman, Huan Wang, and John Wright. “Exact Recovery of Sparsely-Used Dictio-  
406 naries” (June 2012). arXiv: [1206.5882 \[cs.LG\]](https://arxiv.org/abs/1206.5882). 2.
- 407 [18] Joan Bruna and Stéphane Mallat. “Invariant scattering convolution networks”. *IEEE transac-*  
408 *tions on pattern analysis and machine intelligence* 35.8 (Aug. 2013), pp. 1872–1886. 2.
- 409 [19] Peyman Milanfar. “A Tour of Modern Image Filtering: New Insights and Methods, Both  
410 Practical and Theoretical”. *IEEE Signal Processing Magazine* 30.1 (Jan. 2013), pp. 106–128.  
411 4.
- 412 [20] Singanallur V Venkatakrisnan, Charles A Bouman, and Brendt Wohlberg. “Plug-and-Play pri-  
413 ors for model based reconstruction”. *2013 IEEE Global Conference on Signal and Information*  
414 *Processing*. Dec. 2013, pp. 945–948. 4.
- 415 [21] Rémi Gribonval, Rodolphe Jenatton, and Francis Bach. “Sparse and spurious: dictionary  
416 learning with noise and outliers” (July 2014). arXiv: [1407.5155 \[cs.LG\]](https://arxiv.org/abs/1407.5155). 2.
- 417 [22] Jascha Sohl-Dickstein, Eric A Weiss, Niru Maheswaranathan, and Surya Ganguli. “Deep  
418 Unsupervised Learning using Nonequilibrium Thermodynamics” (Mar. 2015). arXiv: [1503.](https://arxiv.org/abs/1503.03585)  
419 [03585 \[cs.LG\]](https://arxiv.org/abs/1503.03585). 1, 4.
- 420 [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. “Deep Residual Learning for  
421 Image Recognition”. *2016 IEEE Conference on Computer Vision and Pattern Recognition*  
422 *(CVPR)*. June 2016, pp. 770–778. 1.

- 423 [24] René Vidal, Yi Ma, and Shankar Sastry. *Generalized Principal Component Analysis*. Springer  
424 Verlag, 2016. 5.
- 425 [25] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. “Mask R-CNN” (Mar. 2017).  
426 arXiv: [1703.06870](https://arxiv.org/abs/1703.06870) [cs.CV]. 1.
- 427 [26] Ilya Loshchilov and Frank Hutter. “Decoupled weight decay regularization”. *arXiv preprint*  
428 *arXiv:1711.05101* (2017). 24.
- 429 [27] Yaniv Romano, Michael Elad, and Peyman Milanfar. “The Little Engine That Could: Regular-  
430 ization by Denoising (RED)”. *SIAM journal on imaging sciences* 10.4 (Jan. 2017), pp. 1804–  
431 1844. 4.
- 432 [28] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
433 Łukasz Kaiser, and Illia Polosukhin. “Attention is all you need”. *Advances in neural informa-*  
434 *tion processing systems* 30 (2017). 1, 3, 6.
- 435 [29] Yubei Chen, Dylan Paiton, and Bruno Olshausen. “The sparse manifold transform”. *Advances*  
436 *in neural information processing systems* 31 (2018). 2.
- 437 [30] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. “Bert: Pre-training of  
438 deep bidirectional transformers for language understanding”. *arXiv preprint arXiv:1810.04805*  
439 (2018). 1.
- 440 [31] Tero Karras, Samuli Laine, and Timo Aila. “A Style-Based Generator Architecture for Genera-  
441 tive Adversarial Networks” (Dec. 2018). arXiv: [1812.04948](https://arxiv.org/abs/1812.04948) [cs.NE]. 1.
- 442 [32] Vardan Papayan, Yaniv Romano, Jeremias Sulam, and Michael Elad. “Theoretical Foundations  
443 of Deep Learning via Sparse Representations: A Multilayer Sparse Model and Its Connection  
444 to Convolutional Neural Networks”. *IEEE Signal Processing Magazine* 35.4 (July 2018),  
445 pp. 72–89. 2.
- 446 [33] Xiaoxia Sun, Nasser M Nasrabadi, and Trac D Tran. “Supervised deep sparse coding networks”.  
447 *2018 25th IEEE International Conference on Image Processing (ICIP)*. IEEE. 2018, pp. 346–  
448 350. 6.
- 449 [34] Yang Song and Stefano Ermon. “Generative Modeling by Estimating Gradients of the Data  
450 Distribution” (July 2019). arXiv: [1907.05600](https://arxiv.org/abs/1907.05600) [cs.LG]. 1.
- 451 [35] John Zarka, Louis Thiry, Tomás Angles, and Stéphane Mallat. “Deep network classification  
452 by scattering and homotopy dictionary learning”. *arXiv preprint arXiv:1910.03561* (2019). 6.
- 453 [36] Lucas Beyer, Olivier J Hénaff, Alexander Kolesnikov, Xiaohua Zhai, and Aäron van den Oord.  
454 “Are we done with imagenet?” *arXiv preprint arXiv:2006.07159* (2020). 24.
- 455 [37] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-  
456 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. “Language  
457 models are few-shot learners”. *Advances in neural information processing systems* 33 (2020),  
458 pp. 1877–1901. 1.
- 459 [38] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and  
460 Sergey Zagoruyko. “End-to-End Object Detection with Transformers” (May 2020). arXiv:  
461 [2005.12872](https://arxiv.org/abs/2005.12872) [cs.CV]. 1.
- 462 [39] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. “A Simple Frame-  
463 work for Contrastive Learning of Visual Representations”. *Proceedings of the 37th Interna-*  
464 *tional Conference on Machine Learning*. Ed. by Hal Daumé Iii and Aarti Singh. Vol. 119.  
465 *Proceedings of Machine Learning Research*. PMLR, 2020, pp. 1597–1607. 1.
- 466 [40] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
467 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,  
468 et al. “An image is worth 16x16 words: Transformers for image recognition at scale”. *arXiv*  
469 *preprint arXiv:2010.11929* (2020). 1, 3, 9.
- 470 [41] Jonathan Ho, Ajay Jain, and Pieter Abbeel. “Denoising diffusion probabilistic models”. *Ad-*  
471 *vances in Neural Information Processing Systems* 33 (2020), pp. 6840–6851. 1.
- 472 [42] Zahra Kadhodaie and Eero P Simoncelli. “Solving Linear Inverse Problems Using the Prior  
473 Implicit in a Denoiser” (July 2020). arXiv: [2007.13640](https://arxiv.org/abs/2007.13640) [cs.CV]. 4.
- 474 [43] Jiaming Song, Chenlin Meng, and Stefano Ermon. “Denoising Diffusion Implicit Models”  
475 (Oct. 2020). arXiv: [2010.02502](https://arxiv.org/abs/2010.02502) [cs.LG]. 1, 4.
- 476 [44] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon,  
477 and Ben Poole. “Score-Based Generative Modeling through Stochastic Differential Equations”  
478 (Nov. 2020). arXiv: [2011.13456](https://arxiv.org/abs/2011.13456) [cs.LG]. 1, 4, 5.

- 479 [45] Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola.  
480 “What makes for good views for contrastive learning?” *Advances in neural information pro-*  
481 *cessing systems* 33 (2020), pp. 6827–6839. 2.
- 482 [46] Yaodong Yu, Kwan Ho Ryan Chan, Chong You, Chaobing Song, and Yi Ma. “Learning Diverse  
483 and Discriminative Representations via the Principle of Maximal Coding Rate Reduction”.  
484 *Advances in Neural Information Processing Systems* 33 (2020), pp. 9422–9434. 2, 3, 5, 17, 22.
- 485 [47] Yuexiang Zhai, Zitong Yang, Zhenyu Liao, John Wright, and Yi Ma. “Complete dictionary  
486 learning via  $l_4$ -norm maximization over the orthogonal group”. *The Journal of Machine*  
487 *Learning Research* 21.1 (2020), pp. 6622–6689. 2.
- 488 [48] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia  
489 Schmid. “Vivit: A video vision transformer”. *Proceedings of the IEEE/CVF international*  
490 *conference on computer vision*. 2021, pp. 6836–6846. 1.
- 491 [49] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. “Masked  
492 Autoencoders Are Scalable Vision Learners” (Nov. 2021). arXiv: 2111.06377 [cs.CV]. 1.
- 493 [50] Geoffrey Hinton. *How to represent part-whole hierarchies in a neural network*. 2021. arXiv:  
494 2102.12627 [cs.CV]. 6.
- 495 [51] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-  
496 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya  
497 Sutskever. “Learning Transferable Visual Models From Natural Language Supervision”. *Pro-*  
498 *ceedings of the 38th International Conference on Machine Learning*. Ed. by Marina Meila and  
499 Tong Zhang. Vol. 139. Proceedings of Machine Learning Research. PMLR, 2021, pp. 8748–  
500 8763. 1.
- 501 [52] Bahareh Tolooshams and Demba Ba. “Stable and Interpretable Unrolled Dictionary Learning”.  
502 *arXiv preprint arXiv:2106.00058* (2021). 2.
- 503 [53] Ilya Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas  
504 Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, Mario Lucic,  
505 and Alexey Dosovitskiy. “MLP-Mixer: An all-MLP Architecture for Vision” (May 2021).  
506 arXiv: 2105.01601 [cs.CV]. 21.
- 507 [54] Kwan Ho Ryan Chan, Yaodong Yu, Chong You, Haozhi Qi, John Wright, and Yi Ma. “ReduNet:  
508 A White-box Deep Network from the Principle of Maximizing Rate Reduction”. *Journal of*  
509 *Machine Learning Research* 23.114 (2022), pp. 1–103. 2, 3, 5, 6, 17, 18.
- 510 [55] Hongrui Chen, Holden Lee, and Jianfeng Lu. “Improved Analysis of Score-based Generative  
511 Modeling: User-Friendly Bounds under Minimal Smoothness Assumptions”. *arXiv preprint*  
512 *arXiv:2211.01916* (2022). 2.
- 513 [56] Yuan Gong, Andrew Rouditchenko, Alexander H Liu, David Harwath, Leonid Karlinsky, Hilde  
514 Kuehne, and James R Glass. “Contrastive audio-visual masked autoencoder”. *The Eleventh*  
515 *International Conference on Learning Representations*. 2022. 1.
- 516 [57] Hangfeng He and Weijie J Su. “A law of data separation in deep learning”. *arXiv preprint*  
517 *arXiv:2210.17020* (2022). 9.
- 518 [58] Geoffrey Hinton. *The Forward-Forward Algorithm: Some Preliminary Investigations*. 2022.  
519 arXiv: 2212.13345 [cs.LG]. 2.
- 520 [59] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. “Elucidating the design space of  
521 diffusion-based generative models”. *arXiv preprint arXiv:2206.00364* (2022). 2, 14.
- 522 [60] Frederic Koehler, Alexander Heckett, and Andrej Risteski. “Statistical Efficiency of Score  
523 Matching: The View from Isoperimetry” (Oct. 2022). arXiv: 2210.00726 [cs.LG]. 2.
- 524 [61] Yi Ma, Doris Tsao, and Heung-Yeung Shum. “On the principles of parsimony and self-  
525 consistency for the emergence of intelligence”. *Frontiers of Information Technology & Elec-*  
526 *tronic Engineering* 23.9 (2022), pp. 1298–1323. 1, 3.
- 527 [62] Druv Pai, Michael Psenka, Chih-Yuan Chiu, Manxi Wu, Edgar Dobriban, and Yi Ma. “Pursuit  
528 of a discriminative representation for multiple subspaces via sequential games”. *arXiv preprint*  
529 *arXiv:2206.09120* (2022). 2.
- 530 [63] Mary Phuong and Marcus Hutter. “Formal algorithms for transformers”. *arXiv preprint*  
531 *arXiv:2207.09238* (2022). 19.
- 532 [64] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer.  
533 “High-resolution image synthesis with latent diffusion models”. *Proceedings of the IEEE/CVF*  
534 *Conference on Computer Vision and Pattern Recognition*. 2022, pp. 10684–10695. 1, 2.

- 535 [65] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed  
536 Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, Tim  
537 Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. “Photorealistic Text-to-Image  
538 Diffusion Models with Deep Language Understanding” (May 2022). arXiv: [2205.11487](https://arxiv.org/abs/2205.11487)  
539 [[cs.CV](#)]. 1.
- 540 [66] Asher Trockman, Devin Willmott, and J Zico Kolter. “Understanding the Covariance Structure  
541 of Convolutional Filters” (Oct. 2022). arXiv: [2210.03651](https://arxiv.org/abs/2210.03651) [[cs.CV](#)]. 21.
- 542 [67] Rene Vidal. *Attention: Self-Expression Is All You Need*. Unpublished; available: <https://openreview.net/forum?id=MmujBCLawFo>. 2022. 1.
- 544 [68] Haoqing Wang, Xun Guo, Zhi-Hong Deng, and Yan Lu. “Rethinking minimal sufficient  
545 representation in contrastive learning”. *Proceedings of the IEEE/CVF Conference on Computer  
546 Vision and Pattern Recognition*. 2022, pp. 16041–16050. 2.
- 547 [69] John Wright and Yi Ma. *High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications*. Cambridge University Press, 2022. 3, 19–21.
- 549 [70] Sitan Chen, Giannis Daras, and Alexandros G Dimakis. “Restoration-Degradation Beyond  
550 Linear Diffusions: A Non-Asymptotic Analysis For DDIM-Type Samplers” (Mar. 2023). arXiv:  
551 [2303.03384](https://arxiv.org/abs/2303.03384) [[cs.LG](#)]. 5.
- 552 [71] Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Yao Liu, Hieu Pham,  
553 Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, et al. “Symbolic discovery of optimization  
554 algorithms”. *arXiv preprint arXiv:2302.06675* (2023). 8, 24.
- 555 [72] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin  
556 Gilmer, Andreas Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al.  
557 “Scaling vision transformers to 22 billion parameters”. *arXiv preprint arXiv:2302.05442*  
558 (2023). 1.
- 559 [73] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson,  
560 Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, Piotr Dollár, and Ross  
561 Girshick. “Segment Anything” (Apr. 2023). arXiv: [2304.02643](https://arxiv.org/abs/2304.02643) [[cs.CV](#)]. 1.
- 562 [74] Hongkang Li, Meng Wang, Sijia Liu, and Pin-Yu Chen. “A Theoretical Understanding of  
563 shallow Vision Transformers: Learning, Generalization, and Sample Complexity”. *arXiv  
564 preprint arXiv:2302.06015* (2023). 1.
- 565 [75] Zonglin Li, Chong You, Srinadh Bhojanapalli, Daliang Li, Ankit Singh Rawat, Sashank J Reddi,  
566 Ke Ye, Felix Chern, Felix Yu, Ruiqi Guo, and Sanjiv Kumar. “The Lazy Neuron Phenomenon:  
567 On Emergence of Activation Sparsity in Transformers”. *The Eleventh International Conference  
568 on Learning Representations*. 2023. 21.
- 569 [76] Ravid Shwartz-Ziv and Yann LeCun. “To Compress or Not to Compress—Self-Supervised  
570 Learning and Information Theory: A Review”. *arXiv preprint arXiv:2304.09355* (2023). 2.
- 571 [77] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. “Consistency models”. *arXiv  
572 preprint arXiv:2303.01469* (2023). 2.

573 **A Technical Details from Section 2**

574 **A.1 Companion to Section 2.2**

575 We first wish to re-iterate the core contributions of our approach in Section 2.2 at a slightly more  
 576 technical level. Connections between denoising and score matching are well-understood [59], and  
 577 computing the optimal denoising function (i.e., the conditional expectation) against a mixture-of-  
 578 Gaussians model is a rather simple computation giving existing tools such as Tweedie’s formula [13].  
 579 These are not our main contributions. Instead, the main contributions of Section 2.2 are two-fold:

- 580 • First, we demonstrate a mechanism to learn representations via denoising within a idealized  
 581 mixture of Gaussian data model for a single token (i.e., with sequence length  $N = 1$ ).
- 582 • Second, we illustrate the similarities between a such-derived representation learning scheme  
 583 and existing self-attention layers within the transformer (with sequence length 1), thus  
 584 demonstrating an interpretation of the self-attention layer as a generalized mechanism to  
 585 denoise against a mixture-of-Gaussian-marginal model for a set of tokens.

586 Now we produce the proofs alluded to in Section 2.2, which mostly form the technical aspects of  
 587 the first listed contribution. To simplify the proofs, we use the following notation correspondences:  
 588  $\mathbf{x} \mapsto \mathbf{z}^\ell$ ,  $\mathbf{z} \mapsto \mathbf{z}^{\ell+1}$ , and  $\sigma \mapsto \sigma^\ell$ .

589 **Proposition 1.** *Let  $\mathbf{u}_1, \dots, \mathbf{u}_K \in \mathbb{R}^d$  be independent and have distribution  $\mathbf{u}_k \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_k)$  for  
 590  $\boldsymbol{\Sigma}_k \succeq \mathbf{0}$ , and let  $\mathbf{z}$  take value  $\mathbf{u}_k$  with probability  $\pi_k > 0$ . Let  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  be independent of  $\mathbf{z}$ .  
 591 Let  $\mathbf{x} \doteq \mathbf{z} + \sigma\mathbf{w}$ . Let  $\mathbf{x} \mapsto q(\mathbf{x})$  be the density of  $\mathbf{x}$ . We define*

$$\mathbf{M}_k \doteq (\boldsymbol{\Sigma}_k + \sigma^2 \mathbf{I}_d)^{-1/2} \quad (18)$$

592 and assume that  $\pi_i \det(\mathbf{M}_i) = \pi_j \det(\mathbf{M}_j)$  for all  $1 \leq i \leq j \leq K$ . Then we have

$$\nabla_{\mathbf{x}} \log q(\mathbf{x}) \quad (19)$$

$$= -[\mathbf{M}_1, \dots, \mathbf{M}_K] \left[ \text{diag} \left( \text{softmax} \left( -\frac{1}{2} \begin{bmatrix} \|\mathbf{M}_1^* \mathbf{x}\|_2^2 \\ \vdots \\ \|\mathbf{M}_K^* \mathbf{x}\|_2^2 \end{bmatrix} \right) \right) \otimes \mathbf{I}_d \right] \begin{bmatrix} \mathbf{M}_1^* \mathbf{x} \\ \vdots \\ \mathbf{M}_K^* \mathbf{x} \end{bmatrix}, \quad (20)$$

593 where  $\otimes$  denotes the Kronecker product, i.e., the block matrix defined by

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} A_{11} \mathbf{B} & \cdots & A_{1n} \mathbf{B} \\ \vdots & \ddots & \vdots \\ A_{m1} \mathbf{B} & \cdots & A_{mn} \mathbf{B} \end{bmatrix} \quad (21)$$

594 *Proof.* Let  $u$  be the multinomial random variable such that  $\mathbf{z} = \mathbf{z}_u$ , so that  $u$  has probability mass  
 595 function  $\pi$ . Then by the law of total probability, we have

$$\nabla_{\mathbf{x}} \log q(\mathbf{x}) = \nabla_{\mathbf{x}} \log \sum_{k=1}^K q(\mathbf{x} | k) \pi_k \quad (22)$$

$$= \frac{\sum_{k=1}^K \pi_k \nabla_{\mathbf{x}} q(\mathbf{x} | k)}{\sum_{k=1}^K q(\mathbf{x} | k) \pi_k} \quad (23)$$

596 where  $q(\mathbf{x} | k)$  is the conditional density of  $\mathbf{x}$  given the event  $\{u = k\}$ . To compute this quantity,  
 597 note that *conditional on the value of  $u$* , we have

$$\mathbf{x} = \mathbf{z}_u + \sigma\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_u + \sigma^2 \mathbf{I}_d). \quad (24)$$

598 Thus we have

$$q(\mathbf{x} | k) = \frac{1}{\sqrt{(2\pi)^d \det(\boldsymbol{\Sigma}_k + \sigma^2 \mathbf{I}_d)}} \exp\left(-\frac{1}{2} \mathbf{x}^* (\boldsymbol{\Sigma}_k + \sigma^2 \mathbf{I}_d)^{-1} \mathbf{x}\right), \quad (25)$$

599 This gives

$$\nabla_{\mathbf{x}} q(\mathbf{x} | k) = -q(\mathbf{x} | k) \cdot (\boldsymbol{\Sigma}_k + \sigma^2 \mathbf{I}_d)^{-1} \mathbf{x}. \quad (26)$$

600 Putting this all together, we get

$$\nabla_{\mathbf{x}} \log q(\mathbf{x}) \quad (27)$$

$$= - \frac{\sum_{k=1}^K q(\mathbf{x} | k) \pi_k \cdot (\boldsymbol{\Sigma}_k + \sigma^2 \mathbf{I}_d)^{-1} \mathbf{x}}{\sum_{k=1}^K q(\mathbf{x} | k) \pi_k} \quad (28)$$

$$= - \frac{\sum_{k=1}^K \pi_k \det(\boldsymbol{\Sigma}_k + \sigma^2 \mathbf{I}_d)^{-1/2} \exp(-\frac{1}{2} \mathbf{x}^* (\boldsymbol{\Sigma}_k + \sigma^2 \mathbf{I}_d)^{-1} \mathbf{x}) \cdot (\boldsymbol{\Sigma}_k + \sigma^2 \mathbf{I}_d)^{-1} \mathbf{x}}{\sum_{k=1}^K \pi_k \det(\boldsymbol{\Sigma}_k + \sigma^2 \mathbf{I}_d)^{-1/2} \exp(-\frac{1}{2} \mathbf{x}^* (\boldsymbol{\Sigma}_k + \sigma^2 \mathbf{I}_d)^{-1} \mathbf{x})}. \quad (29)$$

601 Now define  $\mathbf{M}_k \doteq (\boldsymbol{\Sigma}_k + \sigma^2 \mathbf{I}_d)^{-1/2}$ . With this notation, we have

$$\nabla_{\mathbf{x}} \log q(\mathbf{x}) = - \frac{\sum_{k=1}^K \pi_k \det(\mathbf{M}_k) \exp(-\frac{1}{2} \mathbf{x}^* \mathbf{M}_k \mathbf{M}_k^* \mathbf{x}) \cdot \mathbf{M}_k \mathbf{M}_k^* \mathbf{x}}{\sum_{k=1}^K \pi_k \det(\mathbf{M}_k) \exp(-\frac{1}{2} \mathbf{x}^* \mathbf{M}_k \mathbf{M}_k^* \mathbf{x})} \quad (30)$$

$$= - \frac{\sum_{k=1}^K \pi_k \det(\mathbf{M}_k) \exp(-\frac{1}{2} \|\mathbf{M}_k^* \mathbf{x}\|_2^2) \cdot \mathbf{M}_k \mathbf{M}_k^* \mathbf{x}}{\sum_{k=1}^K \pi_k \det(\mathbf{M}_k) \exp(-\frac{1}{2} \mathbf{x}^* \mathbf{M}_k \mathbf{M}_k^* \mathbf{x})}. \quad (31)$$

602 Given our assumption that each  $\pi_k \det(\mathbf{M}_k)$  is the same, we have

$$\nabla_{\mathbf{x}} \log q(\mathbf{x}) \quad (32)$$

$$= - \frac{\sum_{k=1}^K \pi_k \det(\mathbf{M}_k) \exp(-\frac{1}{2} \|\mathbf{M}_k^* \mathbf{x}\|_2^2) \cdot \mathbf{M}_k \mathbf{M}_k^* \mathbf{x}}{\sum_{k=1}^K \pi_k \det(\mathbf{M}_k) \exp(-\frac{1}{2} \|\mathbf{M}_k^* \mathbf{x}\|_2^2)} \quad (33)$$

$$= - \frac{\sum_{k=1}^K \exp(-\frac{1}{2} \|\mathbf{M}_k^* \mathbf{x}\|_2^2) \cdot \mathbf{M}_k \mathbf{M}_k^* \mathbf{x}}{\sum_{k=1}^K \exp(-\frac{1}{2} \|\mathbf{M}_k^* \mathbf{x}\|_2^2)} \quad (34)$$

$$= - \sum_{k=1}^K \mathbf{e}_k^* \text{softmax} \left( -\frac{1}{2} \begin{bmatrix} \|\mathbf{M}_1^* \mathbf{x}\|_2^2 \\ \vdots \\ \|\mathbf{M}_K^* \mathbf{x}\|_2^2 \end{bmatrix} \right) \mathbf{M}_k \mathbf{M}_k^* \mathbf{x} \quad (35)$$

$$= - [\mathbf{M}_1, \dots, \mathbf{M}_K] \left[ \text{diag} \left( \text{softmax} \left( -\frac{1}{2} \begin{bmatrix} \|\mathbf{M}_1^* \mathbf{x}\|_2^2 \\ \vdots \\ \|\mathbf{M}_K^* \mathbf{x}\|_2^2 \end{bmatrix} \right) \right) \otimes \mathbf{I}_d \right] \begin{bmatrix} \mathbf{M}_1^* \mathbf{x} \\ \vdots \\ \mathbf{M}_K^* \mathbf{x} \end{bmatrix}. \quad (36)$$

603

□

604 Now we provide a final justification for the result cited in Section 2.2.

605 **Approximation 2.** In the setting of Proposition 1, diagonalize  $\boldsymbol{\Sigma}_k = \mathbf{U}_k \boldsymbol{\Lambda}_k \mathbf{U}_k^*$  where  $\mathbf{U}_k \in \mathbb{R}^{d \times p}$   
606 is orthogonal and  $\boldsymbol{\Lambda}_k \succ \mathbf{0} \in \mathbb{R}^{p \times p}$  is diagonal.<sup>9</sup> Then we have the approximation

$$\mathbb{E}[\mathbf{z} | \mathbf{x}] \approx [\mathbf{U}_1, \dots, \mathbf{U}_K] \left[ \text{diag} \left( \text{softmax} \left( \frac{1}{2\sigma^2} \begin{bmatrix} \|\mathbf{U}_1^* \mathbf{x}\|_2^2 \\ \vdots \\ \|\mathbf{U}_K^* \mathbf{x}\|_2^2 \end{bmatrix} \right) \right) \otimes \mathbf{I}_p \right] \begin{bmatrix} \mathbf{U}_1^* \mathbf{x} \\ \vdots \\ \mathbf{U}_K^* \mathbf{x} \end{bmatrix}. \quad (37)$$

607 *Proof.* We have

$$\nabla_{\mathbf{x}} \log q(\mathbf{x}) = - \sum_{k=1}^K \mathbf{e}_k^* \text{softmax} \left( -\frac{1}{2} \begin{bmatrix} \|\mathbf{M}_1^* \mathbf{x}\|_2^2 \\ \vdots \\ \|\mathbf{M}_K^* \mathbf{x}\|_2^2 \end{bmatrix} \right) \mathbf{M}_k \mathbf{M}_k^* \mathbf{x} \quad (38)$$

$$= - \sum_{k=1}^K \mathbf{e}_k^* \text{softmax} \left( -\frac{1}{2\sigma^2} \begin{bmatrix} \|\sigma \mathbf{M}_1^* \mathbf{x}\|_2^2 \\ \vdots \\ \|\sigma \mathbf{M}_K^* \mathbf{x}\|_2^2 \end{bmatrix} \right) \mathbf{M}_k \mathbf{M}_k^* \mathbf{x} \quad (39)$$

$$= - \sum_{k=1}^K \mathbf{e}_k^* \text{softmax} \left( \frac{1}{2\sigma^2} \begin{bmatrix} \|\mathbf{x}\|_2^2 - \|\sigma \mathbf{M}_1^* \mathbf{x}\|_2^2 \\ \vdots \\ \|\mathbf{x}\|_2^2 - \|\sigma \mathbf{M}_K^* \mathbf{x}\|_2^2 \end{bmatrix} \right) \mathbf{M}_k \mathbf{M}_k^* \mathbf{x}. \quad (40)$$

<sup>9</sup>This assumption can be easily relaxed to  $\boldsymbol{\Lambda}_k \succeq \mathbf{0}$  for all  $k$ , but requires some more notation to handle, and the form of the solution does not change. Thus we handle the case where all matrices are full rank for simplicity.

608 Now define  $\mathbf{P}_k \doteq \mathbf{I}_d - \sigma \mathbf{M}_k$ , and let  $\mathbf{U}_k^\perp \in \mathbb{R}^{d \times (d-p)}$  be an orthogonal complement of  $\mathbf{U}_k$ . Then  
609 we have

$$\mathbf{P}_k = \mathbf{I}_d - \sigma \mathbf{M}_k \quad (41)$$

$$= \mathbf{I}_d - \sigma (\boldsymbol{\Sigma}_k + \sigma^2 \mathbf{I}_d)^{-1/2} \quad (42)$$

$$= \mathbf{I}_d - \sigma \left( [\mathbf{U}_k \ \mathbf{U}_k^\perp] \begin{bmatrix} \boldsymbol{\Lambda}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{U}_k^* \\ (\mathbf{U}_k^\perp)^* \end{bmatrix} + \sigma^2 \mathbf{I}_d \right)^{-1/2} \quad (43)$$

$$= \mathbf{I}_d - \sigma \left( [\mathbf{U}_k \ \mathbf{U}_k^\perp] \begin{bmatrix} \boldsymbol{\Lambda}_k + \sigma^2 \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \sigma^2 \mathbf{I}_{d-p} \end{bmatrix} \begin{bmatrix} \mathbf{U}_k^* \\ (\mathbf{U}_k^\perp)^* \end{bmatrix} \right)^{-1/2} \quad (44)$$

$$= \mathbf{I}_d - [\mathbf{U}_k \ \mathbf{U}_k^\perp] \begin{bmatrix} \sigma (\boldsymbol{\Lambda}_k + \sigma^2 \mathbf{I}_p)^{-1/2} & \mathbf{0} \\ \mathbf{0} & \sigma \cdot (\sigma^2)^{-1/2} \mathbf{I}_{d-p} \end{bmatrix} \begin{bmatrix} \mathbf{U}_k^* \\ (\mathbf{U}_k^\perp)^* \end{bmatrix} \quad (45)$$

$$= \mathbf{I}_d - [\mathbf{U}_k \ \mathbf{U}_k^\perp] \begin{bmatrix} (\sigma^{-2} \boldsymbol{\Lambda}_k + \mathbf{I}_p)^{-1/2} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{d-p} \end{bmatrix} \begin{bmatrix} \mathbf{U}_k^* \\ (\mathbf{U}_k^\perp)^* \end{bmatrix} \quad (46)$$

$$= [\mathbf{U}_k \ \mathbf{U}_k^\perp] \begin{bmatrix} \mathbf{I}_p - (\sigma^{-2} \boldsymbol{\Lambda}_k + \mathbf{I}_p)^{-1/2} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{U}_k^* \\ (\mathbf{U}_k^\perp)^* \end{bmatrix} \quad (47)$$

$$\approx [\mathbf{U}_k \ \mathbf{U}_k^\perp] \begin{bmatrix} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{U}_k^* \\ (\mathbf{U}_k^\perp)^* \end{bmatrix} \quad (48)$$

$$= \mathbf{U}_k \mathbf{U}_k^*. \quad (49)$$

610 Thus  $\mathbf{P}_k$  is approximately a projection when  $\sigma$  is small. Under this algebraic relation, we have

$$\nabla_{\mathbf{x}} \log q(\mathbf{x}) \quad (50)$$

$$= - \sum_{k=1}^K \mathbf{e}_k^* \text{softmax} \left( \frac{1}{2\sigma^2} \begin{bmatrix} \|\mathbf{x}\|_2^2 - \|\sigma \mathbf{M}_1^* \mathbf{x}\|_2^2 \\ \vdots \\ \|\mathbf{x}\|_2^2 - \|\sigma \mathbf{M}_K^* \mathbf{x}\|_2^2 \end{bmatrix} \right) \mathbf{M}_k \mathbf{M}_k^* \mathbf{x} \quad (51)$$

$$= - \frac{1}{\sigma^2} \sum_{k=1}^K \mathbf{e}_k^* \text{softmax} \left( \frac{1}{2\sigma^2} \begin{bmatrix} \|\mathbf{x}\|_2^2 - \|(\mathbf{I}_d - \mathbf{P}_1)^* \mathbf{x}\|_2^2 \\ \vdots \\ \|\mathbf{x}\|_2^2 - \|(\mathbf{I}_d - \mathbf{P}_K)^* \mathbf{x}\|_2^2 \end{bmatrix} \right) (\mathbf{I}_d - \mathbf{P}_k) (\mathbf{I}_d - \mathbf{P}_k)^* \mathbf{x} \quad (52)$$

$$\approx - \frac{1}{\sigma^2} \sum_{k=1}^K \mathbf{e}_k^* \text{softmax} \left( \frac{1}{2\sigma^2} \begin{bmatrix} \|\mathbf{P}_1^* \mathbf{x}\|_2^2 \\ \vdots \\ \|\mathbf{P}_K^* \mathbf{x}\|_2^2 \end{bmatrix} \right) (\mathbf{I}_d - \mathbf{P}_k) (\mathbf{I}_d - \mathbf{P}_k)^* \mathbf{x} \quad (53)$$

$$\approx - \frac{1}{\sigma^2} \sum_{k=1}^K \mathbf{e}_k^* \text{softmax} \left( \frac{1}{2\sigma^2} \begin{bmatrix} \|\mathbf{P}_1^* \mathbf{x}\|_2^2 \\ \vdots \\ \|\mathbf{P}_K^* \mathbf{x}\|_2^2 \end{bmatrix} \right) (\mathbf{I}_d - \mathbf{P}_k)^* \mathbf{x} \quad (54)$$

$$= - \frac{\mathbf{x}}{\sigma^2} \sum_{k=1}^K \mathbf{e}_k^* \text{softmax} \left( \frac{1}{2\sigma^2} \begin{bmatrix} \|\mathbf{P}_1^* \mathbf{x}\|_2^2 \\ \vdots \\ \|\mathbf{P}_K^* \mathbf{x}\|_2^2 \end{bmatrix} \right) + \frac{1}{\sigma^2} \sum_{k=1}^K \mathbf{e}_k^* \text{softmax} \left( \frac{1}{2\sigma^2} \begin{bmatrix} \|\mathbf{P}_1^* \mathbf{x}\|_2^2 \\ \vdots \\ \|\mathbf{P}_K^* \mathbf{x}\|_2^2 \end{bmatrix} \right) \mathbf{P}_k^* \mathbf{x} \quad (55)$$

$$= - \frac{1}{\sigma^2} \mathbf{x} + \frac{1}{\sigma^2} \sum_{k=1}^K \mathbf{e}_k^* \text{softmax} \left( \frac{1}{2\sigma^2} \begin{bmatrix} \|\mathbf{P}_1^* \mathbf{x}\|_2^2 \\ \vdots \\ \|\mathbf{P}_K^* \mathbf{x}\|_2^2 \end{bmatrix} \right) \mathbf{P}_k^* \mathbf{x} \quad (56)$$

$$\approx - \frac{1}{\sigma^2} \mathbf{x} + \frac{1}{\sigma^2} \sum_{k=1}^K \mathbf{e}_k^* \text{softmax} \left( \frac{1}{2\sigma^2} \begin{bmatrix} \|\mathbf{U}_1^* \mathbf{x}\|_2^2 \\ \vdots \\ \|\mathbf{U}_K^* \mathbf{x}\|_2^2 \end{bmatrix} \right) \mathbf{U}_k \mathbf{U}_k^* \mathbf{x} \quad (57)$$

$$= - \frac{1}{\sigma^2} \mathbf{x} + \frac{1}{\sigma^2} [\mathbf{U}_1, \dots, \mathbf{U}_K] \left[ \text{diag} \left( \text{softmax} \left( \frac{1}{2\sigma^2} \begin{bmatrix} \|\mathbf{U}_1^* \mathbf{x}\|_2^2 \\ \vdots \\ \|\mathbf{U}_K^* \mathbf{x}\|_2^2 \end{bmatrix} \right) \right) \otimes \mathbf{I}_p \right] \begin{bmatrix} \mathbf{U}_1^* \mathbf{x} \\ \vdots \\ \mathbf{U}_K^* \mathbf{x} \end{bmatrix}. \quad (58)$$

611 Plugging this into Tweedie’s formula, we have

$$\mathbb{E}[\mathbf{z} \mid \mathbf{x}] \approx [\mathbf{U}_1, \dots, \mathbf{U}_K] \left[ \text{diag} \left( \text{softmax} \left( \frac{1}{2\sigma^2} \begin{bmatrix} \|\mathbf{U}_1^* \mathbf{x}\|_2^2 \\ \vdots \\ \|\mathbf{U}_K^* \mathbf{x}\|_2^2 \end{bmatrix} \right) \right) \otimes \mathbf{I}_p \right] \begin{bmatrix} \mathbf{U}_1^* \mathbf{x} \\ \vdots \\ \mathbf{U}_K^* \mathbf{x} \end{bmatrix}. \quad (59)$$

612

□

613 *Remark 3.* Although Approximation 2 is stated as an approximation rather than as a proposition, we  
 614 believe it should be possible without too much extra work to convert it into a statement of asymptotic  
 615 equivalence as  $\sigma \rightarrow 0$  (in particular, holding for  $\sigma$  below the smallest (nonzero) eigenvalue of any  
 616  $\Sigma_k$ . Most approximations taken in the derivation of Approximation 2 can immediately be turned into  
 617 asymptotic claims; the only slightly delicate point is treating the softmax, which can be accomplished  
 618 using standard “high temperature” convergence behavior of the softmax function (in particular, as  
 619  $\sigma \rightarrow 0$  in our expressions, the softmax concentrates on the “best head”).

## 620 A.2 Companion to Section 2.3

621 We again wish to re-iterate the core contribution of our approach in Section 2.3. The application of a  
 622 compression perspective to representation learning has been discussed before, for example in the line  
 623 of maximal coding rate reduction works [46]. In Section 2.3, we provide the following contributions  
 624 and developments to this perspective:

- 625 • We propose a generalized coding rate function  $R^c(\cdot; \mathbf{U}_{[K]})$  which measures the coding rate  
 626 with respect to a set of subspaces  $\mathbf{U}_{[K]}$  as opposed to a set of classes (as in [46, 54]), making  
 627 the underlying formulation unsupervised.
- 628 • We then show how if we adopt the framework of alternating minimization of the sparse rate  
 629 reduction objective, then unrolling the first alternating step — gradient descent on this coding  
 630 rate objective — nearly exactly recovers the common multi-head attention mechanism found  
 631 in transformer networks (except that the query/key/value operators are all the same operation  
 632  $\mathbf{U}_k^*$  now, which we interpret as projection onto a single subspace).

633 In the process of the second contribution, and in the following proofs, we make some simple  
 634 approximations and technical assumptions. The validity of these assumptions may be explored, and  
 635 the approximations refined, altogether providing a more complex (and possibly more performant)  
 636 resulting self-attention like operator. For the sake of technical clarity and simplicity in this work, we  
 637 make perhaps the *simplest possible choices*. As a result, we *do not* claim that our network is optimally  
 638 designed, but rather that the principles we develop in this work (compression, denoising, sparsification,  
 639 unrolled optimization) can provide the backbone for far superior and more interpretable network  
 640 architectures in the future on sundry tasks. As it is, with our straightforward, simple, and interpretable  
 641 design, we still obtain meaningful conceptual results and very solid empirical performance.

642 We now give the derivation of the approximation alluded to in Section 2.3.

643 **Approximation 4.** Let  $\mathbf{Z} \in \mathbb{R}^{d \times N}$  have unit-norm columns, and  $\mathbf{U}_{[K]} = (\mathbf{U}_1, \dots, \mathbf{U}_K)$  such that  
 644 each  $\mathbf{U}_k \in \mathbb{R}^{d \times p}$  is an orthogonal matrix, the  $(\mathbf{U}_k)_{k=1}^K$  are incoherent, and the columns of  $\mathbf{Z}$   
 645 approximately lie on  $\bigcup_{k=1}^K \text{Span}(\mathbf{U}_k)$ . Let  $\gamma = \frac{p}{N\epsilon^2}$ . Let  $\kappa > 0$ . Then

$$\mathbf{Z} - \kappa \nabla_{\mathbf{Z}} R^c(\mathbf{Z} \mid \mathbf{U}_{[K]}) \approx (1 - \kappa\gamma)\mathbf{Z} + \kappa\gamma \text{MSSA}(\mathbf{Z} \mid \mathbf{U}_{[K]}), \quad (60)$$

646 where as in Section 2.3 we have

$$\text{SSA}(\mathbf{Z} \mid \mathbf{U}_k) = (\mathbf{U}_k^* \mathbf{Z}) \text{softmax}((\mathbf{U}_k^* \mathbf{Z})^* (\mathbf{U}_k^* \mathbf{Z})), \quad (61)$$

$$\text{MSSA}(\mathbf{Z} \mid \mathbf{U}_{[K]}) = \gamma [\mathbf{U}_1, \dots, \mathbf{U}_K] \begin{bmatrix} \text{SSA}(\mathbf{Z} \mid \mathbf{U}_1) \\ \vdots \\ \text{SSA}(\mathbf{Z} \mid \mathbf{U}_K) \end{bmatrix}, \quad (62)$$

647 where  $\text{softmax}(\cdot)$  is the softmax operator (applied to each column of an input matrix), i.e.,

$$\text{softmax}(\mathbf{v}) = \frac{1}{\sum_{i=1}^n e^{v_i}} \begin{bmatrix} e^{v_1} \\ \vdots \\ e^{v_n} \end{bmatrix}, \quad (63)$$

$$\text{softmax}([\mathbf{v}_1, \dots, \mathbf{v}_K]) = [\text{softmax}(\mathbf{v}_1), \dots, \text{softmax}(\mathbf{v}_K)]. \quad (64)$$

648 *Proof.* According to (9), the gradient  $\nabla_{\mathbf{Z}} R^c(\mathbf{Z}; \mathbf{U}_{[K]})$  is

$$\nabla_{\mathbf{Z}} R^c(\mathbf{Z}; \mathbf{U}_{[K]}) = \gamma \sum_{k=1}^K \mathbf{U}_k \mathbf{U}_k^* \mathbf{Z} (\mathbf{I} + \gamma (\mathbf{U}_k^* \mathbf{Z})^* (\mathbf{U}_k^* \mathbf{Z}))^{-1}. \quad (65)$$

649 Notice that according to [54], the gradient is precisely the residual of a ridge regression for each  
 650 (projected) token  $\mathbf{U}_k^* \mathbf{z}_i$ ; using other projected tokens  $\mathbf{U}_k^* \mathbf{z}_j$  as the regressors, hence being the residual  
 651 of an auto-regression.

652 However, as we have seen in the work of ReduNet [54], computing the inverse  
 653  $(\mathbf{I} + \gamma (\mathbf{U}_k^* \mathbf{Z})^* (\mathbf{U}_k^* \mathbf{Z}))^{-1}$  can be expensive. Hence for computational efficiency, we may approxi-  
 654 mate it with the first order term of its von Neumann expansion:

$$\nabla_{\mathbf{Z}} R^c(\mathbf{Z}; \mathbf{U}_{[K]}) = \gamma \sum_{k=1}^K \mathbf{U}_k \mathbf{U}_k^* \mathbf{Z} (\mathbf{I} + \gamma (\mathbf{U}_k^* \mathbf{Z})^* (\mathbf{U}_k^* \mathbf{Z}))^{-1} \quad (66)$$

$$\approx \gamma \sum_{k=1}^K \mathbf{U}_k \mathbf{U}_k^* \mathbf{Z} (\mathbf{I} - \gamma (\mathbf{U}_k^* \mathbf{Z})^* (\mathbf{U}_k^* \mathbf{Z})) \quad (67)$$

$$= \gamma \sum_{k=1}^K \mathbf{U}_k (\mathbf{U}_k^* \mathbf{Z} - \gamma \mathbf{U}_k^* \mathbf{Z} [(\mathbf{U}_k^* \mathbf{Z})^* (\mathbf{U}_k^* \mathbf{Z})]) \quad (68)$$

655 Notice that the term  $(\mathbf{U}_k^* \mathbf{Z})^* (\mathbf{U}_k^* \mathbf{Z})$  is the auto-correlation among the projected tokens. As the  
 656 tokens  $\mathbf{Z}$  may be from different subspaces, we would prefer to use only tokens that belong to the  
 657 *same* subspace to regress and compress themselves. Hence we may convert the above correlation  
 658 term into a subspace-membership indicator with a softmax operation, whence (68) becomes

$$\nabla_{\mathbf{Z}} R^c(\mathbf{Z}; \mathbf{U}_{[K]}) \approx \gamma \sum_{k=1}^K \mathbf{U}_k (\mathbf{U}_k^* \mathbf{Z} - \gamma \mathbf{U}_k^* \mathbf{Z} [(\mathbf{U}_k^* \mathbf{Z})^* (\mathbf{U}_k^* \mathbf{Z})]) \quad (69)$$

$$\approx \gamma \sum_{k=1}^K \mathbf{U}_k \mathbf{U}_k^* \mathbf{Z} - \gamma^2 \sum_{k=1}^K \mathbf{U}_k (\mathbf{U}_k^* \mathbf{Z} \text{softmax}((\mathbf{U}_k^* \mathbf{Z})^* (\mathbf{U}_k^* \mathbf{Z}))) \quad (70)$$

659 Then, we can rewrite the above approximation to the gradient of  $R^c$  as:

$$\nabla_{\mathbf{Z}} R^c(\mathbf{Z}; \mathbf{U}_{[K]}) \approx \gamma \sum_{k=1}^K \mathbf{U}_k \mathbf{U}_k^* \mathbf{Z} - \gamma^2 \sum_{k=1}^K \mathbf{U}_k (\mathbf{U}_k^* \mathbf{Z} \text{softmax}((\mathbf{U}_k^* \mathbf{Z})^* (\mathbf{U}_k^* \mathbf{Z}))) \quad (71)$$

$$= \gamma \sum_{k=1}^K \mathbf{U}_k \mathbf{U}_k^* \mathbf{Z} - \gamma^2 \sum_{k=1}^K \mathbf{U}_k \text{SSA}(\mathbf{Z} \mid \mathbf{U}_k) \quad (72)$$

$$= \underbrace{\left( \gamma \sum_{k=1}^K \mathbf{U}_k \mathbf{U}_k^* \right)}_{\approx \gamma \mathbf{Z}} \mathbf{Z} - \gamma^2 [\mathbf{U}_1, \dots, \mathbf{U}_K] \begin{bmatrix} \text{SSA}(\mathbf{Z} \mid \mathbf{U}_1) \\ \vdots \\ \text{SSA}(\mathbf{Z} \mid \mathbf{U}_K) \end{bmatrix} \quad (73)$$

$$\approx \gamma \mathbf{Z} - \gamma^2 [\mathbf{U}_1, \dots, \mathbf{U}_K] \begin{bmatrix} \text{SSA}(\mathbf{Z} \mid \mathbf{U}_1) \\ \vdots \\ \text{SSA}(\mathbf{Z} \mid \mathbf{U}_K) \end{bmatrix}. \quad (74)$$

660 Thus the gradient descent step with learning rate  $\kappa > 0$  gives

$$\mathbf{Z} - \kappa \nabla_{\mathbf{Z}} R^c(\mathbf{Z} \mid \mathbf{U}_{[K]}) \approx (1 - \kappa \gamma) \mathbf{Z} + \kappa \gamma^2 [\mathbf{U}_1, \dots, \mathbf{U}_K] \begin{bmatrix} \text{SSA}(\mathbf{Z} \mid \mathbf{U}_1) \\ \vdots \\ \text{SSA}(\mathbf{Z} \mid \mathbf{U}_K) \end{bmatrix}. \quad (75)$$

661

□

662 **A.3 Companion to Section 2.4**

663 We again wish to re-iterate the core contribution of our approach in Section 2.4.

- 664 • Within the framework of alternating minimization of the sparse rate reduction objective, we
- 665 show that the second alternating step — gradient descent on the overall coding rate plus a
- 666 sparse regularization term — has heuristic connections to a particular LASSO optimization.
- 667 • We show that the unrolling of the proximal gradient step to solve this LASSO optimization
- 668 resembles the MLP which immediately follows the self-attention layer within transformer
- 669 blocks.

670 In the main text, our connection between the second step of the alternating minimization and the  
 671 LASSO optimization was high-level and heuristic. In some sense, the choice to pose the minimization  
 672 step as a LASSO was a *simple, reliable, and interpretable choice* which works well in practice, but  
 673 is nonetheless not backed up by rigorous theoretical justification. In the following subsection, we  
 674 provide a mathematical justification for a reformulation of the minimization step using a majorization-  
 675 minimization framework. We further show that the associated unrolled optimization step bears a  
 676 strong resemblance to the ISTA step. This confirms our earlier discussion — we took the *simplest*  
 677 *possible choice* in designing CRATE, but by more rigorous derivation we can uncover alternative  
 678 operators which nonetheless have the same conceptual function and may perform better in practice.

679 **Assumptions.** In this section, we present a rigorous optimization analysis of an incremental  
 680 minimization approach to the objective (13). We will show that under two simplifying assumptions,  
 681 namely

- 682 1. The columns of  $\mathbf{Z}^{\ell+1/2}$  are normalized, in the sense that  $\text{diag}((\mathbf{Z}^{\ell+1/2})^* \mathbf{Z}^{\ell+1/2}) = \mathbf{1}$ ,<sup>10</sup>
- 683 2. We have  $d \geq N$ ,<sup>11</sup> and the columns of  $\mathbf{Z}^{\ell+1/2}$  are orthogonal, so that  $(\mathbf{Z}^{\ell+1/2})^* \mathbf{Z}^{\ell+1/2} =$   
 684  $\mathbf{I}$ .<sup>12</sup>

685 the approach leads to an update iteration that is equal to a slightly simplified version of the ISTA  
 686 block (16). We see this as a justification for our derivation in Section 2.4, which obtained the ISTA  
 687 block by introducing an additional simplifying assumption on the distribution of the data at layer  $\ell$ .

688 **Analysis.** Following (15), we will consider the natural relaxation of the  $\ell_0$  “norm” to the  $\ell^1$  norm,  
 689 and incorporate a nonnegativity constraint. Consider the objective

$$\varphi(\mathbf{Z}) = \lambda \|\mathbf{Z}\|_1 + \chi_{\{\mathbf{Z} \geq \mathbf{0}\}}(\mathbf{Z}) - \underbrace{\frac{1}{2} \log \det(\mathbf{I} + \alpha \mathbf{Z}^* \mathbf{Z})}_{R(\mathbf{Z})}, \quad (76)$$

690 where  $\mathbf{Z} \in \mathbb{R}^{d \times N}$  and  $\alpha = d/N\varepsilon^2$ , and  $\chi_{\{\mathbf{Z} \geq \mathbf{0}\}}$  denotes the characteristic function for the set of  
 691 elementwise-nonnegative matrices  $\mathbf{Z}$ . As in Appendix A.2, we calculate

$$\nabla_{\mathbf{Z}} R(\mathbf{Z}) = \alpha \mathbf{Z} (\mathbf{I} + \alpha \mathbf{Z}^* \mathbf{Z})^{-1}. \quad (77)$$

692 We consider an incremental optimization scheme for the highly nonlinear and nonconvex objective  $\varphi$ .  
 693 Following Section 2.3, we optimize locally at a “post-compression” iterate  $\mathbf{Z}^{\ell+1/2}$ . We follow the  
 694 standard proximal majorize-minimize framework [69] for incremental/local optimization: this begins  
 695 with the second-order Taylor expansion for the smooth part of  $\varphi$  in a neighborhood of the current

<sup>10</sup>This is a natural assumption in transformer-type architectures such as CRATE due to the use of LayerNorm blocks—although these blocks (indeed, as we use them in CRATE) include trainable mean and scale offsets as well as an additional mean subtraction operation [63], they are initialized to have zero mean and unit norm, hence this assumption corresponds to an analysis of the network at its initialization.

<sup>11</sup>This assumption is without loss of generality, as we will see in the analysis below. The reason is that  $\mathbf{Z}^* \mathbf{Z}$  and  $\mathbf{Z} \mathbf{Z}^*$  have the same nonzero eigenvalues regardless of the shape of  $\mathbf{Z}$ , which implies that  $\log \det(\mathbf{I} + \alpha \mathbf{Z}^* \mathbf{Z}) = \log \det(\mathbf{I} + \alpha \mathbf{Z} \mathbf{Z}^*)$ . In particular, interpreting the norms appropriately (with a slight abuse of notation), we have  $\varphi(\mathbf{Z}) = \varphi(\mathbf{Z}^*)$ , so for the purposes of analysis we can always proceed as though  $\mathbf{Z}$  is a tall matrix (as long as we do not use any special properties of  $\alpha$  in our derivation).

<sup>12</sup>This assumption is strictly stronger than the previous one, and strictly stronger than an assumption of incoherence on the columns. It corresponds to the representation  $\mathbf{Z}^{\ell+1/2}$  being non-collapsed, which we expect to hold at initialization due to the projections  $\mathbf{U}_{[K]}$  being random.

696 iterate  $\mathbf{Z}^{\ell+1/2}$ :

$$R(\mathbf{Z}) = R(\mathbf{Z}^{\ell+1/2}) + \left\langle \nabla_{\mathbf{Z}} R(\mathbf{Z}^{\ell+1/2}), \mathbf{Z} - \mathbf{Z}^{\ell+1/2} \right\rangle + \int_0^1 (1-t) \left\langle \mathbf{Z} - \mathbf{Z}^{\ell+1/2}, \nabla^2 R(\mathbf{Z}_t) (\mathbf{Z} - \mathbf{Z}^{\ell+1/2}) \right\rangle dt, \quad (78)$$

697 where for any  $\mathbf{Z} \in \mathbb{R}^{d \times N}$ ,  $\mathbf{Z}_t = t\mathbf{Z}^{\ell+1/2} + (1-t)\mathbf{Z}$ . The proximal majorization-minimization  
698 approach alternates two steps to minimize  $\varphi$ :

- 699 1. First, use assumptions on  $\mathbf{Z}^{\ell+1/2}$  to derive an upper bound on the operator norm of the  
700 Hessian  $\nabla^2 R(\mathbf{Z})$  over the effective domain of the optimization problem. We will write  $L$   
701 for this (uniform) upper bound. This yields a quadratic upper bound for the smooth part of  
702 the objective  $\varphi$ .
- 703 2. Then, alternately minimize the *smooth part* of the quadratic upper bound as a function of  $\mathbf{Z}$ ,  
704 and take a *proximal step* on the nonsmooth part. It can be shown [69] that corresponds to  
705 the iteration

$$\mathbf{Z}^+ = \text{prox}_{\frac{\lambda}{L}(\|\cdot\|_1 + \chi_{\{\mathbf{Z} \geq \mathbf{0}\}})} \left( \mathbf{Z} + \frac{1}{L} \nabla_{\mathbf{Z}} R(\mathbf{Z}) \right) \quad (79)$$

706 In the alternating minimization setting of this paper for optimizing (1), we only take one  
707 such step, starting at  $\mathbf{Z}^{\ell+1/2}$ .

708 We will instantiate this program below, showing quantitative error bounds related to our assumptions  
709 above as necessary. Rather than directly applying the iteration (79), we will derive it below under our  
710 aforementioned assumptions.

711 Starting at (78), our first task is to upper bound the quadratic residual. This corresponds to estimating

$$\left\langle \mathbf{Z} - \mathbf{Z}^{\ell+1/2}, \nabla^2 R(\mathbf{Z}_t) (\mathbf{Z} - \mathbf{Z}^{\ell+1/2}) \right\rangle \quad (80)$$

$$\leq \sup_{t \in [0,1]} \left\| \nabla^2 R(\mathbf{Z}_t) \right\|_{\ell^2 \rightarrow \ell^2} \left\| \mathbf{Z} - \mathbf{Z}^{\ell+1/2} \right\|_{\mathbb{F}}^2 \quad (81)$$

712 with Cauchy-Schwarz. Using Lemma 5, we can estimate the operator norm term in the previous  
713 bound in terms of properties of  $\mathbf{Z}^{\ell+1/2}$ . We need to bound

$$\alpha \sup_{\|\Delta\|_{\mathbb{F}} \leq 1} \left\| (\Delta - \alpha \mathbf{Z}_t (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} (\mathbf{Z}_t^* \Delta + \Delta^* \mathbf{Z}_t)) (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \right\|_{\mathbb{F}}, \quad (82)$$

714 and Lemma 6 gives that this term is no larger than  $9\alpha/4$  for any  $\mathbf{Z}$  and any  $t$ . With this estimate and  
715 (78), we have a quadratic upper bound for  $-R(\mathbf{Z})$ :

$$-R(\mathbf{Z}) \leq -R(\mathbf{Z}^{\ell+1/2}) + \left\langle -\nabla_{\mathbf{Z}} R(\mathbf{Z}^{\ell+1/2}), \mathbf{Z} - \mathbf{Z}^{\ell+1/2} \right\rangle + \frac{9\alpha}{8} \left\| \mathbf{Z} - \mathbf{Z}^{\ell+1/2} \right\|_{\mathbb{F}}^2. \quad (83)$$

716 Meanwhile, by our assumptions above, we have

$$-\nabla_{\mathbf{Z}} R(\mathbf{Z}^{\ell+1/2}) = -\alpha \mathbf{Z}^{\ell+1/2} (\mathbf{I} + \alpha \mathbf{I})^{-1} = -\frac{\alpha}{1+\alpha} \mathbf{Z}^{\ell+1/2}. \quad (84)$$

717 We now minimize the preceding quadratic upper bound as a function of  $\mathbf{Z}$ . Differentiating, the  
718 minimizer  $\mathbf{Z}_{\text{opt}}$  is calculated as

$$\mathbf{Z}_{\text{opt}} = \left( 1 + \frac{4}{9(1+\alpha)} \right) \mathbf{Z}^{\ell+1/2}, \quad (85)$$

719 and it is well-known that the proximal operator of the sum of  $\chi_{\{\mathbf{Z} \geq \mathbf{0}\}}$  and  $\lambda \|\cdot\|_1$  is simply the  
720 one-sided soft-thresholding operator [69]

$$\text{prox}_{\chi_{\{\mathbf{Z} \geq \mathbf{0}\}} + \lambda \|\cdot\|_1} (\mathbf{Z}) = \max\{\mathbf{Z} - \lambda \mathbf{1}, \mathbf{0}\}, \quad (86)$$

721 where the maximum is applied elementwise. As in Section 2.4, we may write this elementwise  
722 maximum simply as ReLU. Thus, one step of proximal majorization-minimization under our  
723 simplifying assumptions takes the form

$$\mathbf{Z}^{\ell+1} = \text{ReLU} \left( \left( 1 + \frac{4}{9(1+\alpha)} \right) \mathbf{Z}^{\ell+1/2} - \frac{4\lambda}{9\alpha} \mathbf{1} \right). \quad (87)$$

724 Finally, we point out one additional elaboration which introduces the dictionary  $\mathbf{D}$  that appears in the  
725 ISTA block in Section 2.4. Notice that for any orthogonal  $\mathbf{D}$ , one has  $R(\mathbf{D}\mathbf{Z}) = R(\mathbf{Z})$  for every  $\mathbf{Z}$ .  
726 This symmetry implies equivariance properties of  $\nabla_{\mathbf{Z}}R(\mathbf{Z})$  and  $\nabla_{\mathbf{Z}}^2R(\mathbf{Z})$ : for every  $\mathbf{Z}$  and every  $\Delta$   
727 and every orthogonal  $\mathbf{D}$ ,

$$\mathbf{D}\nabla_{\mathbf{Z}}R(\mathbf{Z}) = \nabla_{\mathbf{Z}}R(\mathbf{D}\mathbf{Z}), \quad (88)$$

$$\langle \mathbf{D}\Delta, \nabla_{\mathbf{Z}}^2R(\mathbf{Z})(\mathbf{D}\Delta) \rangle = \langle \Delta, \nabla_{\mathbf{Z}}^2R(\mathbf{D}\mathbf{Z})(\Delta) \rangle. \quad (89)$$

728 Hence the quadratic Taylor expansion (78) can be written equivalently as

$$\begin{aligned} R(\mathbf{Z}) &= R(\mathbf{D}^*\mathbf{Z}^{\ell+1/2}) + \left\langle \nabla_{\mathbf{Z}}R(\mathbf{D}^*\mathbf{Z}^{\ell+1/2}), \mathbf{Z} - \mathbf{Z}^{\ell+1/2} \right\rangle \\ &\quad + \int_0^1 (1-t) \left\langle \mathbf{Z} - \mathbf{Z}^{\ell+1/2}, \nabla^2R(\mathbf{D}^*\mathbf{Z}_t) \left( \mathbf{Z} - \mathbf{Z}^{\ell+1/2} \right) \right\rangle dt, \end{aligned} \quad (90)$$

729 for any orthogonal  $\mathbf{D}$ . The significance of this is that we have obtained an expression equivalent  
730 to (78), but with  $\mathbf{Z}^{\ell+1/2}$  replaced by  $\mathbf{D}^*\mathbf{Z}^{\ell+1/2}$ ; moreover, because our approximation arguments  
731 above are not affected by left-multiplication of  $\mathbf{Z}^{\ell+1/2}$  by an orthogonal matrix (this operation  
732 does not change the norms of the columns of  $\mathbf{Z}^{\ell+1/2}$ , or their correlations, and hence the matrix's  
733 incoherence), we can apply exactly the same line of reasoning above to obtain that an equivalent  
734 proximal majorization-minimization iteration is given by

$$\mathbf{Z}^{\ell+1} = \text{ReLU} \left( \left( 1 + \frac{4}{9(1+\alpha)} \right) \mathbf{D}^*\mathbf{Z}^{\ell+1/2} - \frac{4\lambda}{9\alpha} \mathbf{1} \right), \quad (91)$$

735 for any orthogonal dictionary  $\mathbf{D}$ . This gives an update quite similar to the ISTA block (16) in the  
736 case where the dictionary used in Section 2.4 is orthogonal, but without a skip connection.

737 We thus obtain a natural white-box version of this part of the architecture, along with the natural  
738 interpretation *that its purpose is to sparsify the compressed tokens  $\mathbf{Z}^{\ell+1/2}$  in a (learnable) dictionary*,  
739 which accords with recent empirical studies [75].

740 **Other architectures?** As we mentioned at the start of this section, the preceding derivation  
741 is performed in the most elementary possible setting in order to demonstrate the majorization-  
742 minimization approach for layer design. More precise approximations or assumptions may lead to  
743 superior layer designs that better optimize the target objective (1) (and in particular (13)). We mention  
744 two here:

- 745 1. **Beyond exactly-incoherent features:** our derivations above assumed that the incoming  
746 representations  $\mathbf{Z}^{\ell+1/2}$  were already maximal for the expansion term  $R$  in (13). It is  
747 desirable to obtain a ‘perturbative’ derivation, which applies in cases where  $\mathbf{Z}^{\ell+1/2}$  is not  
748 fully orthogonal, but instead near-orthogonal, in particular *incoherent* [69]. The derivations  
749 above can be adapted to this setting; the perturbation bounds become slightly more delicate,  
750 and the ultimate layer (91) changes to involve additional normalization.
- 751 2. **Beyond orthogonal dictionaries:** The symmetries of the expansion term  $R$  in (13) may be  
752 followed to lead to a pair of dictionaries  $\mathbf{D}$  and  $\mathbf{D}'$  and an objective that sparsifies  $\mathbf{D}\mathbf{Z}\mathbf{D}'$ .  
753 This type of transformation is suggestive of popular architectures that mix over tokens [53,  
754 66], however we consider the simpler form  $\mathbf{D}\mathbf{Z}$  in this work. In addition, we have focused  
755 for simplicity on orthogonal dictionaries  $\mathbf{D}$ ; as in the previous bullet, one may consider  
756 in a similar way dictionaries  $\mathbf{D}$  which are complete and near-orthogonal. Adapting the  
757 derivation to *overcomplete dictionaries* is an interesting future direction that we expect to  
758 improve the scalability of CRATE; one avenue to achieve this could be increasing the number  
759 of projections  $\mathbf{U}_{[K]}$  and their embedding dimensions.

### 760 A.3.1 Auxiliary Lemmas

761 **Lemma 5.** *Consider the function*

$$R(\mathbf{Z}) = \frac{1}{2} \log \det (\mathbf{I} + \alpha \mathbf{Z}^* \mathbf{Z}), \quad (92)$$

762 *where  $\alpha > 0$  is a constant. Then we have*

$$\nabla_{\mathbf{Z}}R(\mathbf{Z}) = \alpha \mathbf{Z} (\mathbf{I} + \alpha \mathbf{Z}^* \mathbf{Z})^{-1}, \quad (93)$$

763 and the Hessian operator  $\nabla_{\mathbf{Z}}^2 R(\mathbf{Z}): \mathbb{R}^{d \times N} \rightarrow \mathbb{R}^{d \times N}$  satisfies that for any  $\Delta \in \mathbb{R}^{d \times N}$ ,

$$\nabla_{\mathbf{Z}}^2 R(\mathbf{Z})(\Delta) \quad (94)$$

$$= \alpha \Delta (\mathbf{I} + \alpha \mathbf{Z}^* \mathbf{Z})^{-1} - \alpha^2 \mathbf{Z} (\mathbf{I} + \alpha \mathbf{Z}^* \mathbf{Z})^{-1} (\mathbf{Z}^* \Delta + \Delta^* \mathbf{Z}) (\mathbf{I} + \alpha \mathbf{Z}^* \mathbf{Z})^{-1}. \quad (95)$$

764 *Proof.* The gradient calculation follows from [46], for example. For the Hessian, we use the usual  
765 approach to calculating derivatives: if  $\Delta$  is any matrix with the same shape as  $\mathbf{Z}$  and  $t > 0$ ,

$$\nabla_{\mathbf{Z}}^2 R(\mathbf{Z})(\Delta) = \left. \frac{\partial}{\partial t} \right|_{t=0} [t \mapsto \nabla_{\mathbf{Z}} R(\mathbf{Z} + t\Delta)], \quad (96)$$

766 valid since  $R$  is smooth. We have

$$\begin{aligned} & \nabla_{\mathbf{Z}} R(\mathbf{Z} + t\Delta) \\ &= \alpha (\mathbf{Z} + t\Delta) (\mathbf{I} + \alpha (\mathbf{Z} + t\Delta)^* (\mathbf{Z} + t\Delta))^{-1} \\ &= \alpha (\mathbf{Z} + t\Delta) (\mathbf{I} + \alpha \mathbf{Z}^* \mathbf{Z} + \alpha t [\mathbf{Z}^* \Delta + \Delta^* \mathbf{Z} + t \Delta^* \Delta])^{-1} \\ &= \alpha (\mathbf{Z} + t\Delta) \left( \mathbf{I} + \alpha t (\mathbf{I} + \alpha \mathbf{Z}^* \mathbf{Z})^{-1} [\mathbf{Z}^* \Delta + \Delta^* \mathbf{Z} + t \Delta^* \Delta] \right)^{-1} (\mathbf{I} + \alpha \mathbf{Z}^* \mathbf{Z})^{-1} \\ &= \alpha (\mathbf{Z} + t\Delta) \left( \sum_{k=0}^{\infty} (-\alpha t)^k \left( (\mathbf{I} + \alpha \mathbf{Z}^* \mathbf{Z})^{-1} [\mathbf{Z}^* \Delta + \Delta^* \mathbf{Z} + t \Delta^* \Delta] \right)^k \right) (\mathbf{I} + \alpha \mathbf{Z}^* \mathbf{Z})^{-1}, \end{aligned}$$

767 where in the fourth line we require that  $t$  is sufficiently close to 0 in order to invoke the Neumann  
768 series. First, notice that the term involving  $\Delta^* \Delta$  does not play a role in the final expression: after  
769 we differentiate with respect to  $t$  and take a limit  $t \rightarrow 0$ , terms arising due to differentiation of  
770  $t \mapsto t \Delta^* \Delta$  go to zero, because whenever the summation index  $k > 0$  we have a term  $(-\alpha t)^k$  that  
771 goes to zero as  $t \rightarrow 0$ . We thus obtain with the product rule

$$\left. \frac{\partial}{\partial t} \right|_{t=0} [t \mapsto \nabla_{\mathbf{Z}} R(\mathbf{Z} + t\Delta)] \quad (97)$$

$$= \alpha \Delta (\mathbf{I} + \alpha \mathbf{Z}^* \mathbf{Z})^{-1} - \alpha^2 \mathbf{Z} (\mathbf{I} + \alpha \mathbf{Z}^* \mathbf{Z})^{-1} (\mathbf{Z}^* \Delta + \Delta^* \mathbf{Z}) (\mathbf{I} + \alpha \mathbf{Z}^* \mathbf{Z})^{-1}. \quad (98)$$

772

□

773 **Lemma 6.** *One has*

$$\sup_{\|\Delta\|_{\mathbb{F}} \leq 1} \left\| (\Delta - \alpha \mathbf{Z}_t (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} (\mathbf{Z}_t^* \Delta + \Delta^* \mathbf{Z}_t)) (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \right\|_{\mathbb{F}} \leq \frac{9}{4}. \quad (99)$$

774 *Proof.* Fix  $\Delta$  satisfying  $\|\Delta\|_{\mathbb{F}} \leq 1$ . By the triangle inequality,

$$\left\| (\Delta - \alpha \mathbf{Z}_t (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} (\mathbf{Z}_t^* \Delta + \Delta^* \mathbf{Z}_t)) (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \right\|_{\mathbb{F}} \quad (100)$$

$$\leq \left\| \Delta (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \right\|_{\mathbb{F}} + \alpha \left\| \mathbf{Z}_t (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} (\mathbf{Z}_t^* \Delta + \Delta^* \mathbf{Z}_t) (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \right\|_{\mathbb{F}}. \quad (101)$$

775 For the first term, we note that

$$\left\| \Delta (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \right\|_{\mathbb{F}} = \left\| ((\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \otimes \mathbf{I}) \text{vec}(\Delta) \right\|_{\mathbb{F}}, \quad (102)$$

776 and since  $(\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \preceq \mathbf{I}$ , we obtain from Cauchy-Schwarz<sup>13</sup>

$$\left\| \Delta (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \right\|_{\mathbb{F}} \leq \|\Delta\|_{\mathbb{F}}. \quad (103)$$

777 We can use a similar idea to control the second term. We have from the triangle inequality

$$\left\| \mathbf{Z}_t (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} (\mathbf{Z}_t^* \Delta + \Delta^* \mathbf{Z}_t) (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \right\|_{\mathbb{F}} \quad (104)$$

$$\leq \left\| \mathbf{Z}_t (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \mathbf{Z}_t^* \Delta (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \right\|_{\mathbb{F}} \quad (105)$$

$$+ \left\| (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \mathbf{Z}_t^* \Delta (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \mathbf{Z}_t^* \right\|_{\mathbb{F}}. \quad (106)$$

<sup>13</sup>Recall that the eigenvalues of a Kronecker product of symmetric matrices are the tensor product of the eigenvalues (with multiplicity).

778 For the first term, we have

$$\| \mathbf{Z}_t (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \mathbf{Z}_t^* \Delta (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \|_{\mathbb{F}} \quad (107)$$

$$= \| ((\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \otimes \mathbf{Z}_t (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \mathbf{Z}_t^*) \text{vec}(\Delta) \|_{\mathbb{F}} \quad (108)$$

$$\leq \sigma_{\max} ((\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1}) \sigma_{\max} (\mathbf{Z}_t (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \mathbf{Z}_t^*) \| \Delta \|_{\mathbb{F}} \quad (109)$$

$$\leq \frac{1}{\alpha} \| \Delta \|_{\mathbb{F}}. \quad (110)$$

779 The last estimate follows from a computation using the SVD of  $\mathbf{Z}_t$ . Meanwhile, we have for the  
 780 second term by a similar argument (using the fact that the singular values of  $\mathbf{A}$  and  $\mathbf{A}^*$  are identical  
 781 for any matrix  $\mathbf{A}$ )

$$\| (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \mathbf{Z}_t^* \Delta (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \mathbf{Z}_t \|_{\mathbb{F}} \leq \sigma_{\max} ((\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \mathbf{Z}_t^*)^2 \| \Delta \|_{\mathbb{F}} \quad (111)$$

$$\leq \frac{1}{4\alpha} \| \Delta \|_{\mathbb{F}}, \quad (112)$$

782 where once again the estimate follows from a computation involving the SVD of  $\mathbf{Z}_t$  (together with  
 783 the fact that the function  $\sigma \mapsto \sigma/(1 + \alpha\sigma^2)$  is bounded on  $\sigma \geq 0$  by  $1/(2\sqrt{\alpha})$ ). Putting it together,  
 784 we have obtained

$$\| (\Delta - \alpha \mathbf{Z}_t (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} (\mathbf{Z}_t^* \Delta + \Delta^* \mathbf{Z}_t)) (\mathbf{I} + \alpha \mathbf{Z}_t^* \mathbf{Z}_t)^{-1} \|_{\mathbb{F}} \leq \frac{9}{4} \| \Delta \|_{\mathbb{F}}, \quad (113)$$

785 which gives the claim after taking suprema.

786

□

## 787 B Additional Experiments and Details

788 In this section, we provide details about our experiments, and report the results of additional experi-  
789 ments that were not covered in the main text. CRATE takes arguably the most basic design choices  
790 possible, and so we do *not* attempt to directly compete with state-of-the-art performance from heavily  
791 engineered and empirically designed transformers. The results of our experiments are meant to  
792 convey a few core messages:

- 793 • *Despite not being engineered to compete with the state-of-the-art, CRATE performs strongly*  
794 *on large-scale real-world datasets, including classification on ImageNet-1K. CRATE also*  
795 *achieves strong transfer learning performance.*
- 796 • *Because our model is designed through unrolled optimization of a well-understood objective,*  
797 *each layer is interpretable.* In particular, we can analyze the performance of CRATE, as well  
798 as design network modifications, on a *layer-wise basis*. This is powered by an arguably  
799 unparalleled level of insight into the role of each operator in our network.
- 800 • *We make the simplest possible choices during the design of CRATE, but these can be changed*  
801 *easily while keeping the same framework.* We study a few modifications later in this section  
802 (Appendix B.4) and show that they do not significantly hurt empirical performance, but  
803 emphasize here that there is significant potential for improvement with different architecture  
804 choices (and in particular a different theoretical analysis).

### 805 B.1 Implementation details

806 In this subsection, we provide more details for implementing CRATE on vision tasks.

#### 807 B.1.1 Architecture of CRATE

808 **Architectural modifications.** Compared to the conceptual architecture proposed in Sections 2.5  
809 and 3, we make the following change for the sake of implementation simplicity:

- 810 • In the compression step, replace the term  $\frac{p}{N\epsilon^2} [\mathbf{U}_1, \dots, \mathbf{U}_K]$  in the MSSA operator with  
811 another trainable parameter  $\mathbf{W} \in \mathbb{R}^{d \times pK}$ . Thus the MSSA block becomes

$$\text{MSSA}(\mathbf{Z} | \mathbf{U}_{[K]}, \mathbf{W}) \doteq \mathbf{W} \begin{bmatrix} \text{SSA}(\mathbf{Z} | \mathbf{U}_1) \\ \vdots \\ \text{SSA}(\mathbf{Z} | \mathbf{U}_K) \end{bmatrix}. \quad (114)$$

812 PyTorch **code for CRATE.** We provide PyTorch-style code for implementing our proposed network  
813 architecture. Algorithm 1 defines the overall architecture, Algorithm 2 and Algorithm 3 contain  
814 details for the transformer block, self-attention block (MSSA-block), and MLP block (ISTA-block).

#### 815 B.1.2 Training Setup

816 **Pre-training on ImageNet-1K.** We apply the Lion optimizer [71] for pre-training both CRATE and  
817 ViT models. We configure the learning rate as  $2.4 \times 10^{-4}$ , weight decay as 0.5, and batch size as  
818 2,048. We incorporate a warm-up strategy with a linear increase over 5 epochs, followed by training  
819 the models for a total of 150 epochs with cosine decay. For data augmentation, we only apply the  
820 standard techniques, random cropping and random horizontal flipping, on the ImageNet-1K dataset.  
821 We apply label smoothing with smoothing parameter 0.1. One training epoch of CRATE-Base takes  
822 around 240 seconds using 16 A100 40GB GPUs.

823 **Fine-tuning.** We fine-tune our pre-trained CRATE and ViT models on the following target datasets:  
824 CIFAR10/CIFAR100 [10], Oxford Flowers-102 [7], Oxford-IIIT-Pets [16]. We also evaluate our  
825 pre-trained models on the commonly used ImageNet Real [36] benchmark. For each fine-tuning  
826 task, we use the AdamW optimizer [26]. We configure the learning rate as  $5 \times 10^{-5}$ , weight decay  
827 as 0.01, and batch size to be 512. To allow transfer learning, we first resize our input data to  
828 224. For data augmentations, we also adopt several standard techniques: random cropping, random  
829 horizontal flipping, and random augmentation (with number of transformations  $n = 2$  and magnitude  
830 of transformations  $m = 14$ ).<sup>14</sup>

<sup>14</sup>[https://github.com/huggingface/pytorch-image-models/blob/main/timm/data/auto\\_ augment.py](https://github.com/huggingface/pytorch-image-models/blob/main/timm/data/auto_augment.py)

---

**Algorithm 1:** PyTorch-style pseudocode for CRATENetwork

---

```
# Class ViT_dictionary definition
CRATE:
# initialization
def init(self, image_size, patch_size, num_classes, dim, depth, heads,
    mlp_dim, pool = 'cls', channels = 3, dim_head = 64, dropout = 0.,
    emb_dropout = 0.):
    # define patch, image dimensions and number of patches
    image_height, image_width = pair(image_size)
    patch_height, patch_width = pair(patch_size)
    num_patches = (image_height // patch_height) * (image_width //
        patch_width)
    patch_dim = channels * patch_height * patch_width

    # define patch embedding, positional embedding, dropout, and transformer
    self.to_patch_embedding = Sequential(Rearrange, LayerNorm(patch_dim),
        Linear(patch_dim, dim), LayerNorm(dim))
    self.pos_embedding = Parameter(random(1, num_patches + 1, dim))
    self.cls_token = Parameter(random(1, 1, dim))
    self.dropout = Dropout(emb_dropout)
    self.transformer = Transformer(dim, depth, heads, dim_head, mlp_dim,
        dropout)

    # define pooling, latent layer, and MLP head
    self.pool = pool
    self.to_latent = Identity()
    self.mlp_head = Sequential(LayerNorm(dim), Linear(dim, num_classes))

# forward pass
def forward(self, img):
    x = self.to_patch_embedding(img)
    b, n, _ = shape(x)
    cls_tokens = repeat(self.cls_token, '1 1 d -> b 1 d', b = b)
    x = concatenate((cls_tokens, x), dim=1)
    x += self.pos_embedding[:, :(n + 1)]
    x = self.dropout(x)
    x = self.transformer(x)
    x = mean(x, dim = 1) if self.pool == 'mean' else x[:, 0]
    x = self.to_latent(x)
    return self.mlp_head(x)
```

---

---

**Algorithm 2:** Pytorch Style Pseudocode for Transformer Block in CRATE

---

```
# Class Transformer definition
class Transformer:
    # initialization
    def init(self, dim, depth, heads, dim_head, mlp_dim, dropout = 0.):
        # define layers
        self.layers = []
        self.depth = depth
        for _ in range(depth):
            self.layers.append([LayerNorm(dim, Attention(dim, heads, dim_head,
                dropout))])
            self.layers.append([LayerNorm(dim, FeedForward(dim, mlp_dim,
                dropout))])

    # forward pass
    def forward(self, x):
        for attn, ff in self.layers:
            x_ = attn(x) + x
            x = ff(x_)
        return x
```

---

---

**Algorithm 3:** Pseudocode for Attention and FeedForward

---

```
# Class FeedForward definition
class FeedForward:
    # initialization
    def init(self, dim, hidden_dim, dropout = 0., step_size=0.1, lambd=0.1):
        self.weight = Parameter(Tensor(dim, dim))
        init.kaiming_uniform_(self.weight)
        self.step_size = step_size
        self.lambd = lambd
    # forward pass
    def forward(self, x):
        x1 = linear(x, self.weight, bias=None)
        grad_1 = linear(x1, self.weight.t(), bias=None)
        grad_2 = linear(x, self.weight.t(), bias=None)
        grad_update = self.step_size * (grad_2 - grad_1) - self.step_size *
            self.lambd
        output = relu(x + grad_update)
        return output
# Class Attention definition
class Attention:
    # initialization
    def init(self, dim, heads = 8, dim_head = 64, dropout = 0.):
        inner_dim = dim_head * heads
        project_out = not (heads == 1 and dim_head == dim)
        self.heads = heads
        self.scale = dim_head ** -0.5
        self.attend = Softmax(dim = -1)
        self.dropout = Dropout(dropout)
        self.qkv = Linear(dim, inner_dim, bias=False)
        self.to_out = Sequential(Linear(inner_dim, dim), Dropout(dropout)) if
            project_out else nn.Identity()
    # forward pass
    def forward(self, x):
        w = rearrange(self.qkv(x), 'b n (h d) -> b h n d', h = self.heads)
        dots = matmul(w, w.transpose(-1, -2)) * self.scale
        attn = self.attend(dots)
        attn = self.dropout(attn)
        out = matmul(attn, w)
        out = rearrange(out, 'b h n d -> b n (h d)')
        return self.to_out(out)
```

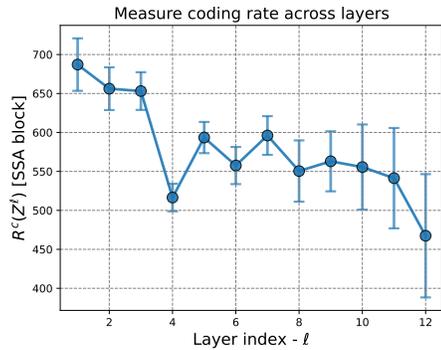
---

831 **B.2 Experimental Results**

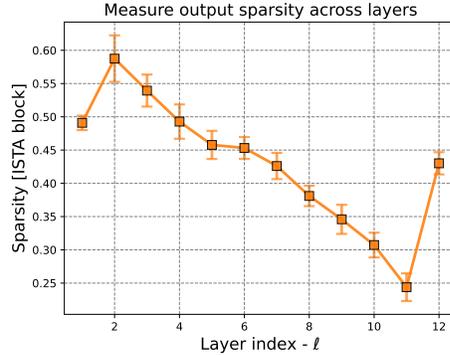
832 In this subsection, we provide additional experimental results on CRATE, including layer-wise  
 833 measurements, visualizations, as well as ablation studies.

834 **B.2.1 Layer-wise Evaluation and Visualization**

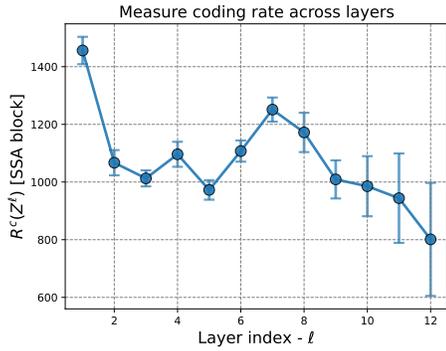
835 **Layer-wise evaluation of compression and sparsity.** Similar to Figure 3, we conduct the layer-  
 836 wise evaluation of compression term and sparsity for CRATE-Tiny, CRATE-Base, and CRATE-Large.  
 837 We observe similar behavior as mentioned in Section 3.1: both the compression term and the sparsity  
 838 term improves as the layer index increases.



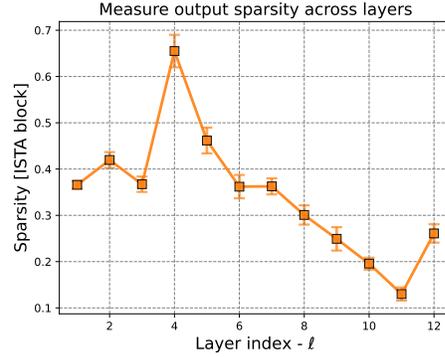
(a) Compression (Model: CRATE-Tiny).



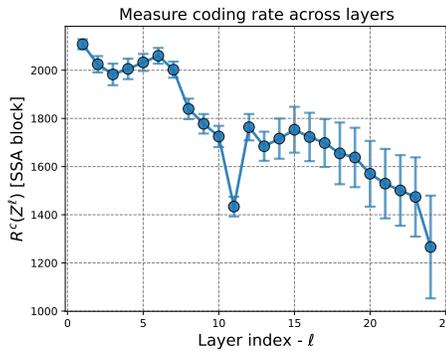
(b) Sparsity (Model: CRATE-Tiny).



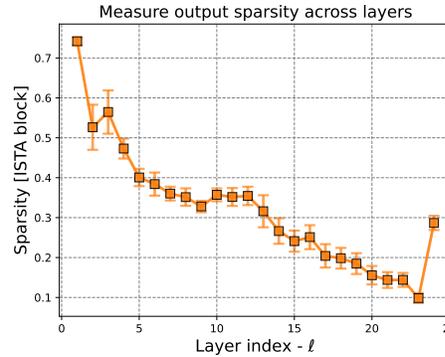
(c) Compression (Model: CRATE-Base).



(d) Sparsity (Model: CRATE-Base).



(e) Compression (Model: CRATE-Large).

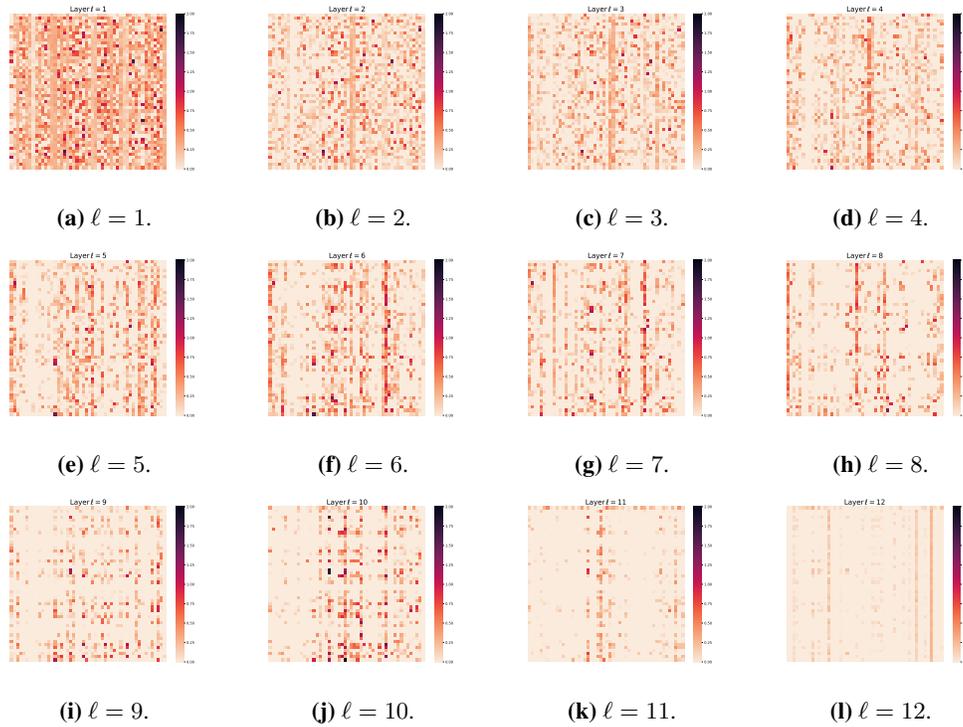


(f) Sparsity (Model: CRATE-Large).

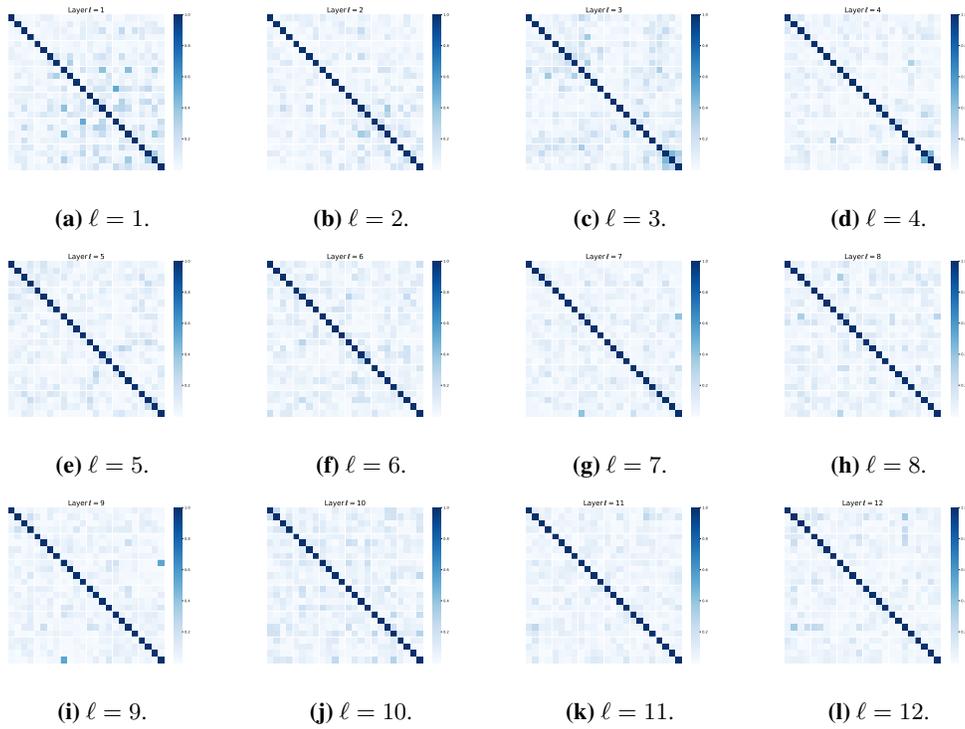
**Figure 5:** *Left:* The compression term  $R^c(\mathbf{Z}^{\ell+1/2})$  of the MSSA outputs at different layers. *Right:* the sparsity of the ISTA output block,  $\|\mathbf{Z}^{\ell+1}\|_0/(d \cdot N)$ , at different layers.

839 **Visualizing layer-wise token representations.** In Figure 6, we visualize the token representations  
 840  $Z^\ell$  at different layers  $\ell \in \{1, \dots, 12\}$ . We provide more results evaluated on other samples in  
 841 Appendix B.2.2.

842 **Visualizing layer-wise subspaces in multi-head self-attention.** We provide the visualization of  
 843  $U_{[K]}^\ell$  in Figure 7.



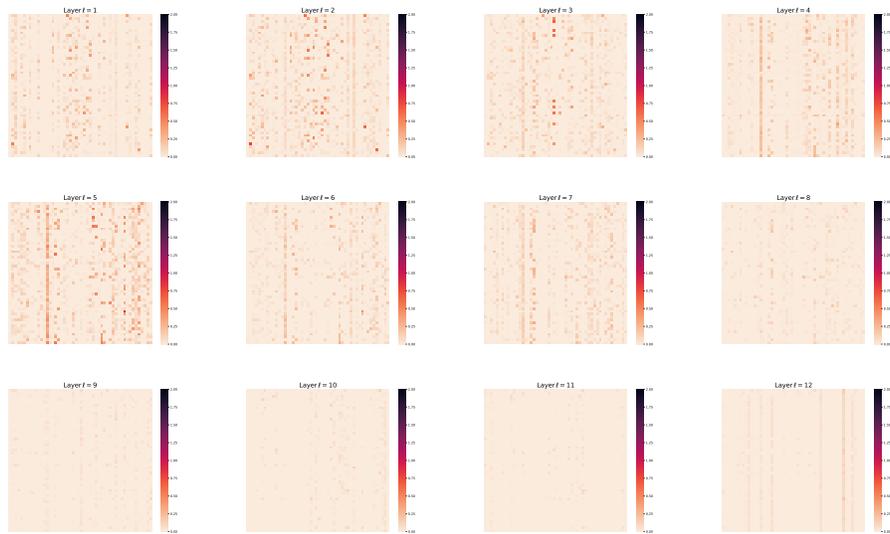
**Figure 6:** Visualizing layer-wise token  $Z^\ell$  representations at each layer  $\ell$ . To enhance the visual clarity, we randomly extract a  $50 \times 50$  sub-matrix from  $Z^\ell$  for display purposes. (Model: CRATE-Tiny)



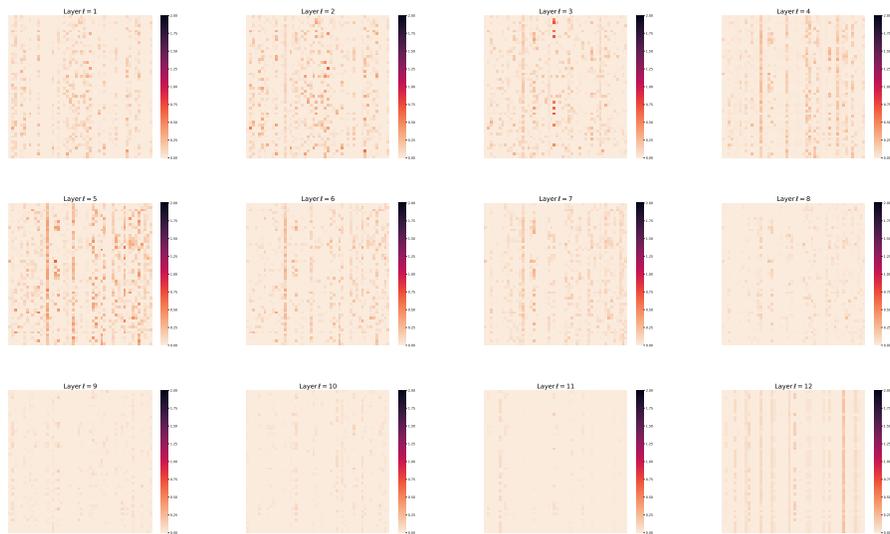
**Figure 7:** We visualize the  $[\mathbf{U}_1^\ell, \dots, \mathbf{U}_K^\ell] [\mathbf{U}_1^\ell, \dots, \mathbf{U}_K^\ell]^* \in \mathbb{R}^{pK \times pK}$  at different layers. The  $(i, j)$ -th block in each sub-figure corresponds to  $(\mathbf{U}_i^\ell)^* \mathbf{U}_j^\ell$  for  $i, j \in [K]$  at a particular layer  $\ell$ . To enhance the visual clarity, for each subspace  $\mathbf{U}_i$ , we randomly pick 4 directions for display purposes. (Model: CRATE-Tiny)

844 **B.2.2 Additional Layer-wise Visualization**

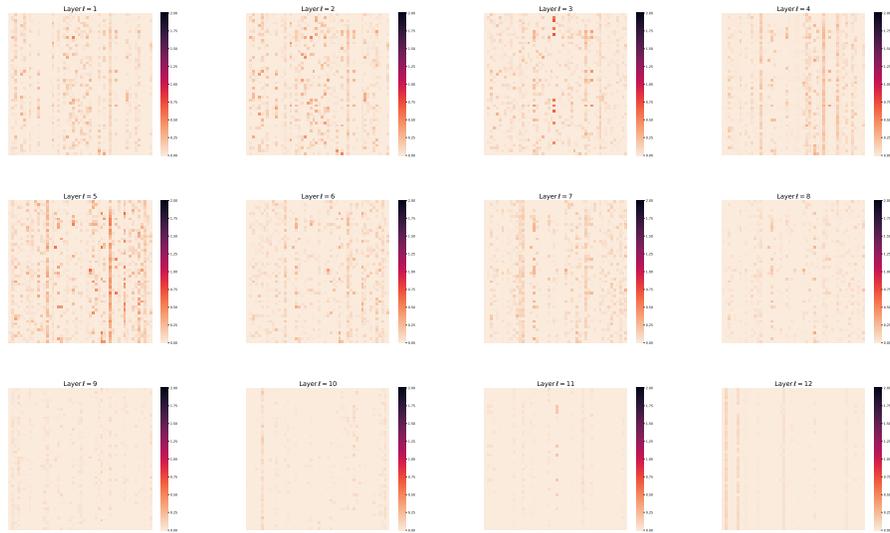
845 We provide more results of the layer-wise token representation visualization on other samples in  
 846 Figure 8, Figure 9, Figure 10, and Figure 11 (Model: CRATE-Base).



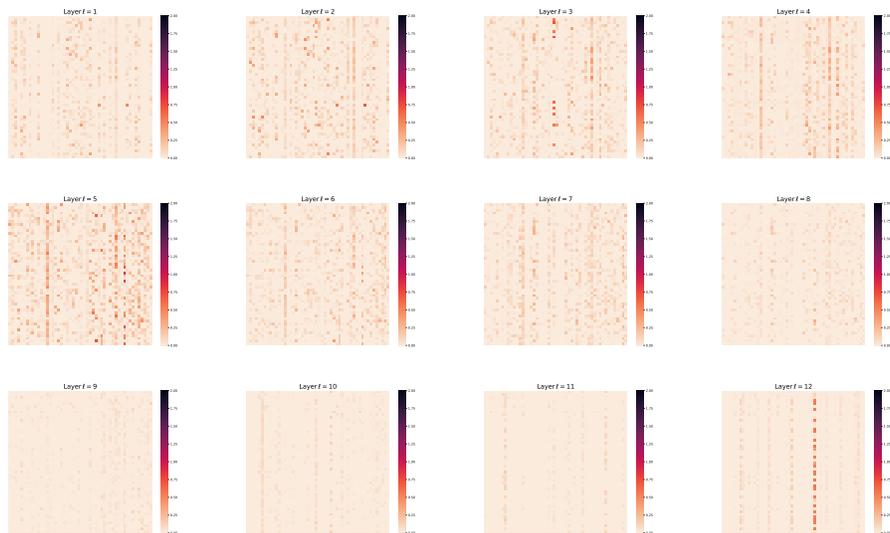
**Figure 8:** Visualizing layer-wise token  $Z^\ell$  representations at each layer  $\ell$ . To enhance the visual clarity, we randomly extract a  $50 \times 50$  sub-matrix from  $Z^\ell$  for display purposes. (*Sample 1*)



**Figure 9:** Visualizing layer-wise token  $Z^\ell$  representations at each layer  $\ell$ . To enhance the visual clarity, we randomly extract a  $50 \times 50$  sub-matrix from  $Z^\ell$  for display purposes. (*Sample 2*)



**Figure 10:** Visualizing layer-wise token  $Z^\ell$  representations at each layer  $\ell$ . To enhance the visual clarity, we randomly extract a  $50 \times 50$  sub-matrix from  $Z^\ell$  for display purposes. (*Sample 3*)



**Figure 11:** Visualizing layer-wise token  $Z^\ell$  representations at each layer  $\ell$ . To enhance the visual clarity, we randomly extract a  $50 \times 50$  sub-matrix from  $Z^\ell$  for display purposes. (*Sample 4*)

847 **B.3 CRATE Ablation**

848 **Hyperparameters of CRATE.** In Table 2, we present evaluation of CRATE trained with various  
 849 parameters. More specifically, we investigate the effect of number of epochs, weight decay, learning  
 850 rate, step size ( $\eta$ ) and the regularization term ( $\lambda$ ) in ISTA block. As shown in Table 2, CRATE  
 851 demonstrates consistently satisfactory performance across a diverse range of hyperparameters.

**Table 2:** Top 1 accuracy of CRATE on various datasets with different architecture design variants when trained on ImageNet.

Model	epoch	weight decay	lr	$\eta$ (ISTA)	$\lambda$ (ISTA)	ImageNet
CRATE-B	150 (default)	0.5 (default)	$2.4 \times 10^{-4}$	0.1	0.1	70.8
CRATE-B	150	0.5	$2.4 \times 10^{-4}$	0.02	0.1	70.7
CRATE-B	150	0.5	$2.4 \times 10^{-4}$	0.5	0.1	66.7
CRATE-B	150	0.5	$2.4 \times 10^{-4}$	0.1	0.02	70.8
CRATE-B	150	0.5	$2.4 \times 10^{-4}$	0.1	0.5	70.5
CRATE-B	90	0.5	$2.4 \times 10^{-4}$	0.1	0.1	69.5
CRATE-B	300	0.5	$2.4 \times 10^{-4}$	0.1	0.1	70.9
CRATE-B	150	1.0	$2.4 \times 10^{-4}$	0.1	0.1	70.3
CRATE-B	150	0.05	$2.4 \times 10^{-4}$	0.1	0.1	70.2
CRATE-B	150	0.5	$4.8 \times 10^{-4}$	0.1	0.1	70.2
CRATE-B	150	0.5	$1.2 \times 10^{-4}$	0.1	0.1	70.3

852 **B.4 Exploring Architecture Variants**

853 In this section, we explore the two following alternative architectures. One architecture involves a  
 854 modification to the attention mechanism, while the other involves a modification to the sparsification  
 855 mechanism. Again, we re-emphasize that these choices, although principled, are entirely modular and  
 856 the choices we make here still lead to very simple architectures. A more sophisticated analysis may  
 857 lead to different, more complicated architectures that perform better in practice. The architectures we  
 858 experiment with are:

- 859 • Compression-inspired attention mechanism: revert the change in (114). That is, the attention  
 860 mechanism implements (11) and (12) directly.
- 861 • Majorization-minimization proximal step sparsification: instead of (16), implement (91).

862 We obtain the following classification results in Table 3. After conducting additional simplifications  
 863 to the network architecture (i.e., imposing additional constraints to the network architecture design),  
 864 we discover that CRATE maintains reasonable performance on ImageNet-1K.

**Table 3:** Top 1 accuracy of CRATE on various datasets with different architecture design variants when trained on ImageNet.

Model	MSSA-block	ISTA-block	ImageNet
CRATE-B	default	default	70.8
CRATE-B	Eq. (11) and (12)	default	63.3
CRATE-B	default	Eq. (91)	68.6