

LEARNING SINGLE-COMPONENT DISCRIMINATIVE REPRESENTATIONS VIA MAXIMAL CODING RATE REDUCTION

Anonymous authors

Paper under double-blind review

ABSTRACT

The recently proposed maximal coding rate reduction principle (MCR²) offers a promising theoretical framework for interpreting modern deep networks through the lens of data compression and discriminative representation. It maps high-dimensional multi-class data into mutually orthogonal linear subspaces, with each subspace capturing as many structural details of its class as possible. In this work, we show that such structural maximization not only increases model sensitivity to feature noise but also hinders generalization. In contrast, we argue that retaining only the single most discriminative structural component per class improves both generalization and robustness to feature noise, while preserving the desirable properties of MCR², such as robustness to label noise and resistance to catastrophic forgetting. We formalize this approach as a new framework termed SiMCoding and validate it extensively across supervised learning, white-box architectures, and incremental learning on diverse datasets. The superior performance of SiMCoding highlights its potential as a strong alternative for medium-scale classification tasks, particularly under label and feature noise.

1 INTRODUCTION

Numerous research efforts have sought to demystify the black-box nature of deep learning. Among these, an influential direction is the principle of *Maximal Coding Rate Reduction* (MCR²) (Yu et al., 2020), which reformulates the learning objective to explicitly capture the low-dimensional structures underlying high-dimensional data, rather than focusing primarily on label fitting. MCR² is grounded in the manifold hypothesis, which posits that although data points $\mathbf{x} \in \mathbb{R}^D$ are observed in a high-dimensional ambient space, their variability is largely confined to a union of low-dimensional submanifolds, $\mathcal{M} = \bigcup_{i=1}^K \mathcal{M}_i$, as illustrated in Figure 1 (Hein & Audibert, 2005; Spigler et al., 2019; Pope et al., 2021; Wright & Ma, 2021). Each submanifold \mathcal{M}_i corresponds to a semantic class or cluster, and the central objective of MCR² is to faithfully uncover and effectively organize these structures in the feature space.

As a foundational concept in information theory, the lossy *coding rate* $R(\mathbf{Z}, \epsilon)$ quantifies **the volume of a distribution** or its finite set \mathbf{Z} , up to a precision ϵ (Rissanen, 1998; Cover, 1999; Ma et al., 2007): a lower coding rate indicates a more compact set. What distinguishes the coding rate from other classical concepts in information theory, such as entropy and mutual information, is that it serves as a well-defined measure of distribution compactness even for degenerate distributions, which commonly arise in data with relatively low intrinsic dimensionality. Formally, in a K -class classification task, let the features of the i -th class be denoted by $\mathbf{Z}_i \in \mathbb{R}^{d \times m_i}$, and define the overall feature set as $\mathbf{Z} = \bigcup_{i=1}^K \mathbf{Z}_i$. The MCR² framework aims to **maximize the volume** of the overall feature set \mathbf{Z} while simultaneously **minimizing the volumes** of the individual cluster sets \mathbf{Z}_i . This simple mechanism effectively maps high-dimensional data into a compact and structured low-dimensional representation, as depicted in Figure 1:

1. **Discriminative representation:** Features of each class \mathbf{Z}_i are compressed into a low-dimensional linear subspace \mathcal{S}_i , and these subspaces are mutually orthogonal, i.e., $\mathbf{Z}_i \mathbf{Z}_j^\top = \mathbf{0}$ for all $i \neq j$.
2. **Diverse representation:** The dimensionality (or variance) of features of each class is maximized subject to the constraint of the representation space \mathbb{R}^d , i.e., $\sum_{i=1}^K \text{rank}(\mathbf{Z}_i) = d$.

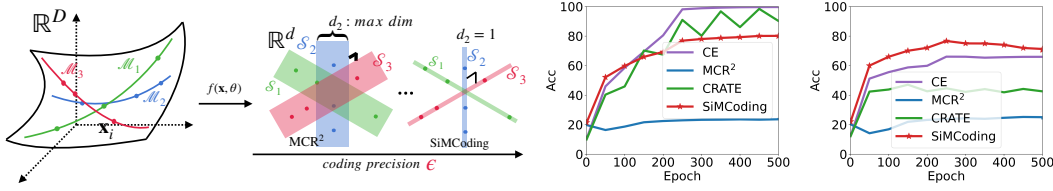


Figure 1: (Left) MCR^2 maps data x_i , typically distributed over nonlinear low-dimensional submanifolds \mathcal{M}_i , onto mutually orthogonal linear subspaces \mathcal{S}_i with maximal dimensionality, whereas SiMCoding enforces each subspace to be one-dimensional. (Middle and Right) Training and test accuracy on CIFAR-20 with 20% randomly corrupted labels.

Note that, to achieve the second property, the features must be encoded with high precision ϵ , enabling the model to capture as many structural details as possible and thereby allowing each class-specific feature set \mathcal{Z}_i to attain its maximal dimensionality, i.e., maximal structural components.

Owing to its simplicity and conceptual interpretability, MCR^2 has emerged as an influential framework in representation learning and has been applied across diverse settings. It has inspired the design of interpretable white-box network architectures (Chan et al., 2022; Pai et al., 2023; Yu et al., 2024a; Yang et al., 2024) and efficient self-attention modules (Wu et al., 2024). It has also been explored in incremental learning (Wu et al., 2021; Tong et al., 2023), generative modelling (Dai et al., 2022), and unsupervised learning (Tong et al., 2022; Chu et al., 2024; Wu et al., 2025), among others.

However, in this work, we question whether it is truly necessary for a model to maximize structural details across different learning settings. In unsupervised learning, where labels are unavailable, it is natural for the model to retain as much structural information as possible in order to deeply uncover the underlying structure and subsequently cluster the samples. In contrast, in supervised learning, where labels provide guidance, preserving only a single discriminative structural component may be sufficient for accurate classification. As a thought experiment, consider classifying images of the digits 0 and 1. Recognizing their overall outlines is sufficient for the task, whereas fine-grained structural details—such as the precise curvature of a 0 or the thickness of a 1—are unnecessary. Even though in practice the features learned by neural networks are often highly abstract and difficult to interpret directly (Zeiler & Fergus, 2014; Chen et al., 2023), this example illustrates that in supervised classification tasks, high-level features may only need to preserve a single discriminative structural component rather than all structural details of the input.

The main contributions of this work are summarized as follows:

- While MCR^2 has achieved remarkable success, particularly in unsupervised learning, we find that in supervised classification its pursuit of maximal structural detail leads to severe underfitting and poor generalization. As shown in Figure 1, MCR^2 struggles to fit CIFAR-20 dataset (with 20% randomly corrupted labels) (Krizhevsky et al., 2009) effectively and exhibits weak generalization. Moreover, it is highly vulnerable to input noise (Table 2).

- To address these issues, we propose learning only the single most discriminative structural component for each class, rather than maximizing all structural details. We term this approach **Single-component Maximal Coding** rate reduction (SiMCoding).

Specifically, we first provide a theoretical analysis showing how the coding precision ϵ determines the extent to which the model emphasizes structural details, formally corresponding to the varying dimensionality of each class subspace \mathcal{S}_i . This analysis further reveals that ϵ can be pre-specified to ensure that each class attains its minimal one dimensional subspace. As an important byproduct, this theory removes the need to tune ϵ as a hyperparameter, thereby significantly reducing the burden of applying SiMCoding.

- We validate SiMCoding across a wide range of datasets and learning settings. Experiments show that SiMCoding matches the fitting ability and generalization of the widely used cross-entropy framework, while exhibiting substantially stronger robustness to label noise. As shown in Figure 1, on CIFAR-20 with 20% randomly corrupted labels, the training accuracy of SiMCoding plateaus near 80%, indicating that it fits only the correctly labeled samples. Moreover, despite adopting the opposite strategy of retaining only a single structural com-

ponent per class, SiMCoding preserves key properties of MCR², including robustness to catastrophic forgetting in incremental learning setting.

- We analyse the computational complexity of SiMCoding and conclude that, despite potential limitations, it remains a strong alternative for classification on datasets with a moderate number of classes (e.g., $K \leq 100$), particularly in the presence of feature or label noise.

2 METHOD

Representation learning and the MCR² principle. We are given data $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m] \in \mathbb{R}^{D \times m}$ consisting of m samples from K classes. The aim of deep representation learning is to transform high-dimensional data into low-dimensional features that capture intrinsic properties such as structure and geometry to facilitate downstream tasks such as classification. A widely used viewpoint, often referred to as the *manifold hypothesis*, suggests that each class lies on a low-dimensional submanifold \mathcal{M}_i , and that the entire dataset is concentrated near the union $\mathcal{M} = \cup_{i=1}^K \mathcal{M}_i$ (Hein & Audibert, 2005; Spigler et al., 2019; Pope et al., 2021; Wright & Ma, 2021). This motivates seeking features $\mathbf{z}_i \in \mathbb{R}^d$ with $d \ll D$ that retain this structure while discarding redundant variability. To obtain such features, one typically employs a nonlinear map $f_{\Theta} : \mathbb{R}^D \rightarrow \mathbb{R}^d$ parametrized by neural network weights Θ :

$$\mathbf{x} \mapsto \mathbf{z} = f_{\Theta}(\mathbf{x}),$$

and collects the feature matrix $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_m] \in \mathbb{R}^{d \times m}$. Desirable features should not only align with the underlying class structure but also admit a compact, structured and interpretable form.

The principle of *Maximal Coding Rate Reduction* (MCR²) (Yu et al., 2020; Chan et al., 2022) provides an information-theoretic criterion for achieving this goal. It simultaneously maximizes the overall volume of all features to encourage separation across classes (*expansion*) and minimizes the average volume of each class to promote compactness (*compression*). Specifically, let $\mathbf{\Pi} = \{\mathbf{\Pi}_i \in \mathbb{R}^{m \times m}\}_{i=1}^K$ denote a set of diagonal matrices, where each diagonal entry $\mathbf{\Pi}_i(j, j)$ specifies the probability that sample j belongs to class i . The MCR² framework then seeks to optimise

$$\begin{aligned} \max_{\mathbf{Z}, \mathbf{\Pi}} \Delta \mathcal{R}(\mathbf{Z}, \mathbf{\Pi}, \epsilon) = & \underbrace{\frac{1}{2} \log \det \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z} \mathbf{Z}^\top \right)}_{\text{Expansion: } \mathcal{R}(\mathbf{Z}, \epsilon)} \\ & - \underbrace{\sum_{i=1}^K \frac{\text{tr}(\mathbf{\Pi}_i)}{2m} \log \det \left(\mathbf{I} + \frac{d}{\text{tr}(\mathbf{\Pi}_i)\epsilon^2} \mathbf{Z} \mathbf{\Pi}_i \mathbf{Z}^\top \right)}_{\text{Compression: } \mathcal{R}_c(\mathbf{Z}, \mathbf{\Pi}, \epsilon)}. \end{aligned} \quad (1)$$

In this formulation, the membership matrices $\mathbf{\Pi}$ may either be fixed by labels (supervised case) or optimised jointly with \mathbf{Z} (unsupervised case). This flexibility enables MCR² to unify both paradigms within a single framework. The coding precision $\epsilon > 0$ is typically and heuristically chosen to be very small so that all fine structural details of data are preserved in learned features.

Coding rate. A central component of MCR² is the *coding rate* $\mathcal{R}(\cdot, \epsilon)$, which quantifies the effective volume of a distribution or its finite sample set under a prescribed distortion level $\epsilon > 0$ (Rissanen, 1998; Cover, 1999; Ma et al., 2007). Formally, for each $\mathbf{z} \in \mathbf{Z}$, let its reconstruction $\hat{\mathbf{z}}$ satisfy

$$\mathbb{E}[\|\mathbf{z} - \hat{\mathbf{z}}\|] \leq \epsilon,$$

the average number of binary bits required to encode the feature set \mathbf{Z} is given by $\mathcal{R}(\mathbf{Z}, \epsilon) = \frac{1}{2} \log \det \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z} \mathbf{Z}^\top \right)$. This expression admits a clear geometric interpretation: it represents the volume of the subspace spanned by \mathbf{Z} , measured in units of ϵ -balls (i.e., d -dimensional spheres of radius ϵ). Intuitively, a larger coding rate indicates that more ϵ -balls are required to cover the feature subspace, implying a richer feature set. This closed-form formulation, originally derived for Gaussian data supported on a subspace (Ma et al., 2007), offers both computational tractability and geometric as well as statistical interpretability within the MCR² framework.

Normalization and geometric view. The coding rate is closely related to the *volume* spanned by the features. If the features are arbitrarily scaled, the measured volumes are no longer comparable across

classes. To ensure fairness, Yu et al. (2020) normalise the scale of each class such that $\|\mathbf{Z}_i\|_F^2 = m_i$, a condition that can be conveniently enforced using batch normalization during training.

From this perspective, the two terms in equation 1 provide a natural geometric interpretation. The expansion term $\mathcal{R}(\mathbf{Z}, \epsilon)$ measures the overall volume of the feature set \mathbf{Z} ; maximizing it encourages the features to spread out and occupy as large a region of the space \mathbb{R}^d as possible. The compression term $\mathcal{R}_c(\mathbf{Z}, \mathbf{\Pi}, \epsilon)$ measures the volume of features within each class; minimizing it reduces the within-class spread, pulling samples of the same label into a compact, low-dimensional cluster.

Yu et al. (2020); Chan et al. (2022) showed that optimizing the overall objective yields representations with two key properties: (i) features within each class concentrate on a linear subspace that reflects the underlying submanifold, while subspaces of different classes tend toward orthogonality, thereby enhancing discriminability; and (ii) with sufficiently high coding precision (e.g., $\epsilon^2 = 0.5$), the subspaces collectively expand to span the full dimensionality of the feature space \mathbb{R}^d , i.e.,

$$\text{rank}(\mathbf{Z}) = \sum_{i=1}^K \text{rank}(\mathbf{Z}_i) = \sum_{i=1}^K d_i = d,$$

so that each class preserves the maximal possible structural components in its feature set \mathbf{Z}_i .

However, we demonstrate that in the supervised setting, blindly maximizing structural details within MCR² can make the model overly sensitive to input noise (Table 2) and may even lead to severe underfitting and poor generalization (Figure 1 and Table 1). To mitigate these issues, we argue that it is sufficient for MCR² to capture only the most discriminative structural component of each class, i.e., $\text{rank}(\mathbf{Z}_i) = 1$. This view resonates with the information bottleneck theory (Tishby & Zaslavsky, 2015; Hu et al., 2024), which states that the role of a neural network is to extract features \mathbf{Z} that retain only the minimal sufficient information relevant to the target labels while discarding irrelevant details. In a similar spirit, but from a different perspective, the MCR² framework aims to capture discriminative low-dimensional structures; thus, in the supervised case, it may be sufficient to preserve only the minimal structural information required for class separation.

To learn single-component discriminative representations via MCR², we present Theorem 1 to characterize how the coding precision ϵ influences the dimensionality of each subspace:

Theorem 1. Let $\mathbf{Z}^* = \mathbf{Z}_1^* \cup \dots \cup \mathbf{Z}_K^*$ be the optimal solution to equation 1. Define $d_i^* = \sqrt{\frac{m_i}{m}} \frac{d}{\epsilon^2}$. Then the following properties hold:

- **Discriminativeness:** Features from different classes reside in mutually orthogonal, low-dimensional linear subspaces; that is, $(\mathbf{Z}_i^*)^\top \mathbf{Z}_j^* = \mathbf{0}$ for all $i \neq j$.
- **Bounded Dimensionality:** Each class-specific subspace has dimensionality $d_i \leq d_i^*$. Furthermore, for each class, the first $d_i - 1$ singular values of \mathbf{Z}_i^* are identical.

This theorem implies that the dimensionality d_i of each class-specific subspace cannot exceed d_i^* . Recall that $\|\mathbf{Z}_i\|_F^2 = m_i = \sum_{j=1}^{\min(d, m_i)} \sigma_j^2$ where σ_j denotes the j -th singular value of \mathbf{Z}_i . As d_i^* decreases, $\text{rank}(\mathbf{Z}_i)$ also decreases, causing the energy to concentrate on fewer singular values. This, in turn, highlights the significance of the remaining structural components. Now we can encourage each class to collapse toward its minimal one-dimensional subspace, i.e., $d_i = \text{rank}(\mathbf{Z}_i) \rightarrow 1$. It should be emphasized that when the data are imbalanced, a uniform coding precision ϵ cannot simultaneously enforce all subspaces to be one-dimensional, since the dimensionality of each class subspace is also influenced by its proportion in the dataset, i.e., $\frac{m_i}{m}$. Therefore, we require the weaker condition $\min_{1 \leq i \leq K} d_i \geq 1$. From this, an upper bound for ϵ can be readily obtained:

$$\epsilon^2 \leq \epsilon_U^2 = \min_{1 \leq i \leq K} d \sqrt{\frac{m_i}{m}}.$$

Theoretically, the feature set \mathbf{Z}_i of each class can be constrained to retain only a single discriminative structural component by setting $\epsilon = \epsilon_U$. However, such a constraint may be overly restrictive in practice, especially for complex deep neural networks with nonlinear objectives, where perfect convergence is rarely attainable on challenging datasets. For instance, Zhu et al. (2021) showed that the global optimality conditions for cross-entropy with certain regularization terms require the existence of at least one redundant dimension in the representation space, along which the gradient can escape local minima. Motivated by this analysis, we advocate setting the minimum dimensionality of

each subspace to $d_i \geq 2$. We emphasize that although Wang et al. (2024) provided a global landscape analysis for a variant of MCR², their formulation differs from ours, and their optimality condition relies on maximising structural details, offering limited guidance for our setting. Moreover, since d_i^* may not serve as a strict upper bound except in the trivial case $d_i^* = d_i = 1$, ensuring $d_i = 2$ requires $d_i^* > 2$. Consequently, we recommend adopting $d_i^* = 3$, which provides both a stronger theoretical guarantee and greater practical robustness.

Accordingly, the practically upper bound for ϵ is

$$\epsilon^2 \leq \epsilon_\star^2 = \min_{1 \leq i \leq K} \frac{d}{3} \sqrt{\frac{m_i}{m}}.$$

As demonstrated in our experiments in Section 3, setting $\epsilon = \epsilon_\star$ is sufficient for the model to learn single-component discriminative representations in practice. Further increasing ϵ does not provide additional benefit.

Building on above theoretical insights, we propose a paradigm shift from conventional MCR², which seeks to capture maximal structural components, toward focusing on the most discriminative component. We refer to this variant as *Single-component Maximal Coding rate reduction* (SiMCoding):

$$\max_{\mathbf{Z}, \mathbf{\Pi}} \Delta \mathcal{R}(\mathbf{Z}, \mathbf{\Pi}, \epsilon) = \mathcal{R}(\mathbf{Z}, \epsilon) - \mathcal{R}_c(\mathbf{Z}, \mathbf{\Pi}, \epsilon), \quad \text{s.t.} \quad \epsilon^2 = \min_{1 \leq i \leq K} \frac{d}{3} \sqrt{\frac{m_i}{m}}. \quad (2)$$

Positioning SiMCoding among MCR² and CRATE. MCR² can be directly employed as a loss function to train predefined neural networks such as ResNet-18 (He et al., 2016), which is arguably the most straightforward way to utilize it. Beyond this, Chan et al. (2022) demonstrated that a deep neural network can be interpreted as the unrolling of iterative gradient steps for optimizing MCR², where each layer corresponds to one iteration. This perspective enables the principled design of *white-box neural networks*. In particular, when shift-invariance is enforced for classification, the resulting architecture naturally takes the form of a multi-channel convolutional network, termed *ReduNet*. A key limitation, however, is that constructing ReduNet is computationally demanding, which hinders its scalability. Consequently, Chan et al. (2022) primarily introduced ReduNet as a rigorous proof-of-concept, with validation limited to small-scale datasets such as MNIST (LeCun, 1998). Nonetheless, this work has been highly influential.

Building on this line of research, Yu et al. (2023; 2024a) introduced sparse MCR² and approximation techniques, giving rise to a white-box transformer-like architecture termed *CRATE*, which has since been applied across diverse domains (Pai et al., 2023; Yu et al., 2024b; Yang et al., 2024). However, CRATE requires learning class-specific sets of orthonormal bases \mathbf{U}_k , which in practice demands large-scale datasets for effective training. In CRATE, the features are projected into these low-dimensional basis spaces \mathbf{U}_k . Since \mathbf{Z} is sparse, the class-specific projection $\mathbf{U}_k^\top \mathbf{Z}$ has lower dimensionality than \mathbf{U}_k itself. Interestingly, in our experiments we observed that CRATE also tends to learn nearly one-dimensional, mutually orthogonal subspaces, albeit through a mechanism fundamentally different from that of SiMCoding.

Computational Complexity. The computational complexity of SiMCoding is dominated by the computation of $(K + 1)$ log-determinants, resulting in a total cost of $\mathcal{O}(K \min(d^3, m^3))$. This indicates that SiMCoding is most suitable for datasets with a moderate number of classes.

It can be concluded that, in the supervised setting, MCR² remains constrained to relatively simple datasets such as MNIST, whereas CRATE represents a significant step toward scalability on large-scale datasets such as ImageNet-1K and ImageNet-21K (Deng et al., 2009). Our proposed SiMCoding strikes a balance, being particularly well-suited for datasets with a moderate number of classes (e.g., $K \leq 100$), where it achieves superior overall performance compared not only to MCR² and CRATE but also to cross-entropy, especially in the presence of label or feature noise.

3 EXPERIMENT

In this section, we evaluate SiMCoding in terms of fitting ability, generalization, and robustness to feature noise, label noise, and catastrophic forgetting. Our aim is not to exhaustively explore extensions or engineering refinements, but rather to demonstrate that even the simplest use of SiMCoding

Table 1: Training and test accuracy (%) of different methods across datasets. Best results are in bold.

Method	MNIST		CIFAR-10		CIFAR-20		CIFAR-100		ImageNette	
	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test
CE	100.00	99.08	99.98	93.23	99.95	80.76	99.98	75.14	97.99	91.78
MCR ²	99.81	98.41	93.83	90.83	43.66	41.89	6.15	6.27	56.72	51.15
CRATE	99.70	96.01	93.88	79.90	93.01	53.45	98.36	51.23	92.17	82.19
SiMCoding	99.94	99.09	99.32	93.04	96.15	80.87	98.90	74.10	95.70	92.85

provides a strong alternative for moderate-scale classification tasks. We compare against influential baselines—cross-entropy (CE), MCR², and CRATE—focusing on validating the effectiveness of the SiMCoding principle under fair conditions. Additional implementation details and experiments are provided in the Appendix, with code included in the supplementary material to reproduce all results.

Performance Metric. Traditional MCR²-based methods typically rely on a Nearest Subspace Classifier (NSC), which assigns labels by measuring the distance of a feature representation to the principal subspace of each class (Yu et al., 2020; Chan et al., 2022). For class i , let μ_i be the class mean and U_i the matrix containing its top r_i principal components. Given a test feature $f(\mathbf{x}'; \theta)$, classification is performed by finding the subspace that minimizes the projection error:

$$i' = \arg \min_i \left\| (I - U_i U_i^\top) (f(\mathbf{x}'; \theta) - \mu_i) \right\|_2^2.$$

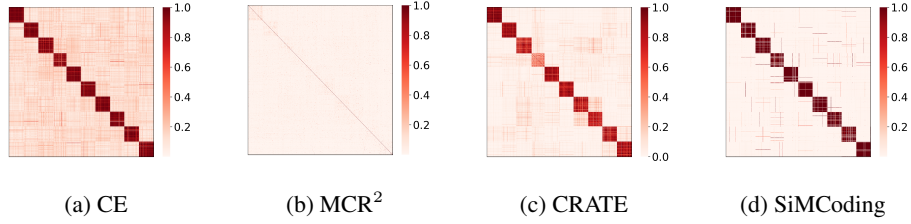
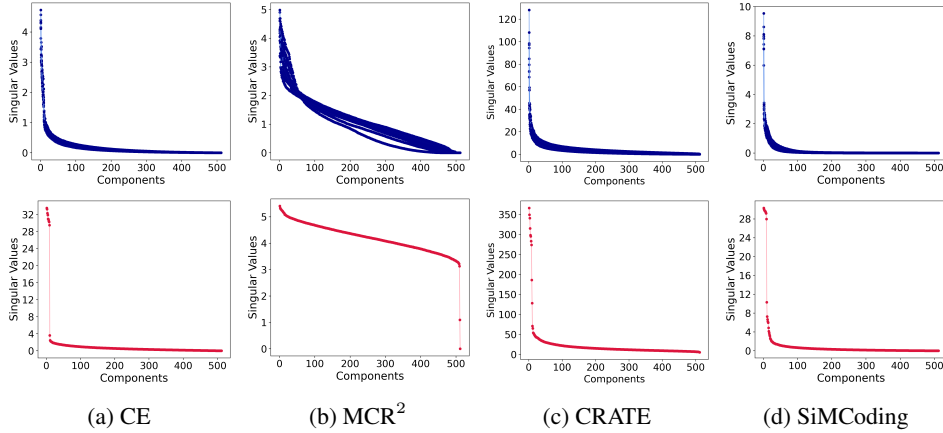
In contrast, for our SiMCoding, the dimensionality d_i of each class subspace is designed to approach one, making it non-trivial to predefine a fixed r_i . As a result, the NSC is not applicable in this setting. Instead, we evaluate both SiMCoding and MCR² by training a simple logistic softmax classifier on their learned features and reporting its accuracy as the performance measure.

3.1 FITTING AND GENERALISATION

Dataset. In this subsection, we study the fitting and generalization ability of SiMCoding. We evaluate on MNIST (LeCun, 1998), the CIFAR family (Krizhevsky et al., 2009) including CIFAR-10, CIFAR-20 (the coarse-label version of CIFAR-100), and CIFAR-100, as well as ImageNette (Howard, 2019), a 10-class subset of ImageNet.

Architecture and Training. For the MNIST dataset, we use a compact convolutional network consisting of two 3×3 convolutional layers (with 32 and 64 channels, both with ReLU activation), followed by 2×2 max pooling, a fully connected ReLU layer with 1024 units, and a final projection to $d = 64$ dimensions. In all experiments, the learned features are normalized such that $\|\mathbf{Z}_i\|_F^2 = m_i$ for each class. For networks trained with CE, we append a classification layer to the same backbone architecture used for MCR² and SiMCoding. For the CIFAR datasets and ImageNette, we employ a ResNet-18 (He et al., 2016) backbone, replacing the final layer with a two-layer ReLU-activated MLP that outputs 512-dimensional representations. To ensure consistency and comparability, we adopt training hyperparameters closely following (Yu et al., 2020; Chan et al., 2022; Yu et al., 2024a). Implementation details are in the Appendix.

Performance Comparison. Table 1 report the final training and test accuracy, and Figure 5 in the Appendix illustrates learning dynamics across datasets. For training accuracy, CE achieves nearly perfect fitting, while SiMCoding attains a comparable level even on challenging datasets such as CIFAR-100 and ImageNette, indicating a fitting capacity on par with CE. CRATE also shows strong fitting but is weaker than CE and SiMCoding, whereas MCR² converges to much lower values, especially on CIFAR-20 and CIFAR-100. In terms of test accuracy, SiMCoding consistently ranks among the best. It matches CE on MNIST and CIFAR-10, while slightly surpassing it on CIFAR-20 and ImageNette, suggesting stronger robustness as dataset complexity increases. CRATE underperforms on test accuracy (e.g., 51.23% on CIFAR-100), reflecting its reliance on class-specific orthonormal bases U_i , which work better with large-scale data. MCR² yields the weakest generalization, consistent with its limited training performance. Overall, SiMCoding combines strong fitting capacity with consistently high generalization, outperforming or matching CE and clearly surpassing CRATE and MCR² in both stability and robustness.

Figure 2: Heatmaps of $|\mathbf{Z}\mathbf{Z}^\top|$ on ImageNette, with samples sorted by class.Figure 3: Feature spectra analysis on ImageNette under different training objectives. Top row: per-class singular value spectra of \mathbf{Z}_i . Bottom row: overall singular value spectrum of \mathbf{Z} .

Structure analysis of LDR. To examine the structure of the learned representations, we compute $\mathbf{Z}\mathbf{Z}^\top$ on ImageNette, sorting samples by class index (Figure 2). This highlights inter-class orthogonality. SiMCoding produces features that are clearly orthogonal across classes, yielding well-separated representations. CRATE also induces a block-diagonal structure, though less sharply, consistent with its slightly lower training accuracy. In contrast, MCR² fails to enforce separation, showing little block structure, while CE exhibits approximate orthogonality, consistent with neural collapse phenomenon (Papayan et al., 2020). Figure 3 shows the singular value spectra of per-class features \mathbf{Z}_i and the overall matrix \mathbf{Z} . For CE, CRATE, and SiMCoding, each class spectrum is dominated by a single leading singular value, indicating nearly rank-one features, while MCR² retains a much flatter spectrum, suggesting high-rank intra-class variability. At the global level, CE, CRATE, and SiMCoding exhibit about ten dominant singular values, aligning with the number of classes and indicating an effectively low-rank feature space, whereas MCR² produces a full-rank spectrum, reflecting its failure to compress intra-class variation or separate classes.

Results on other datasets showing similar patterns are in Figures 7, 8, 9, and 10 in the Appendix.

Discarding Structural details To further investigate the learning mechanism of SiMCoding, we visualize mean saliency maps for 500 randomly selected images per class. For CE, MCR², and SiMCoding, saliency maps are computed directly. For CRATE, which relies on tokenization, we instead aggregate the four attention heads on MNIST into a single visualization. Figure 4 shows the results for digits 0 and 1, with the full set in Figure 6 in the Appendix. The differences are evident. CE produces saliency patterns that only partially align with digit structure, reflecting its focus on label fitting rather than structural abstraction. MCR² captures fine-grained details. CRATE distributes attention across digit components, capturing complementary cues. In contrast, SiMCoding concentrates on the most discriminative element—the digit outline—suggesting a mechanism that filters redundant details while preserving class-defining features.

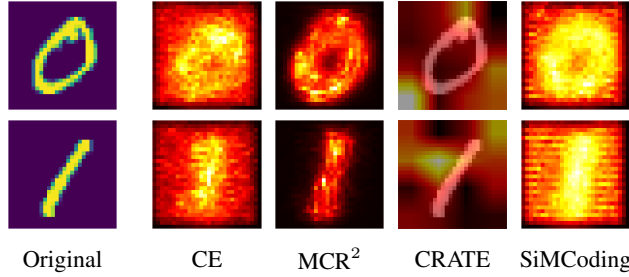


Figure 4: Comparison of saliency maps for digits 0 and 1.

Table 2: Training and test accuracy (%) on CIFAR-10 under different levels of feature noise (std).

Noise Std	Training Accuracy					Test Accuracy				
	0.04	0.08	0.12	0.16	0.20	0.04	0.08	0.12	0.16	0.20
CE	99.91	99.74	99.38	98.60	97.57	91.04	87.94	84.59	81.04	78.59
MCR ²	91.31	87.54	83.79	80.38	75.91	88.30	83.93	79.92	75.29	69.76
CRATE	92.92	90.48	87.46	83.48	79.53	77.64	71.77	62.66	54.62	52.81
SiMCoding	98.26	96.15	93.30	90.16	87.25	91.41	88.02	84.66	81.00	78.03

3.2 ROBUSTNESS AGAINST FEATURE NOISE

Following [Chan et al. \(2022\)](#), we corrupt CIFAR-10 with additive Gaussian noise $\mathcal{N}(0, \sigma^2)$ at $\sigma \in \{0.04, 0.08, 0.12, 0.16, 0.20\}$, keeping the architecture and training setup unchanged. Final accuracies are in Table 2, with learning dynamics in Figure 11 in the Appendix.

It is clear that test accuracy decreases monotonically with σ for all frameworks. CE maintains nearly saturated training accuracy (e.g., 99.9% \rightarrow 97.6%), reflecting its tendency to fit even heavily corrupted inputs. SiMCoding achieves the strongest or tied generalization at low–moderate noise ($\sigma \leq 0.12$) and stays within 0.5 pp of CE at higher noise.

MCR² and CRATE degrade sharply. MCR² encodes noise along with fine details (88.30 \rightarrow 69.76), while CRATE is even more brittle (77.64 \rightarrow 52.81), reflecting its reliance on class bases U_i that also capture corrupted variability. Their convergence is slower and less stable than CE and SiMCoding. Overall, CE and SiMCoding are the most robust, with SiMCoding matching or surpassing CE at moderate noise and maintaining a substantially smaller generalization gap at higher corruption levels.

3.3 ROBUSTNESS AGAINST LABEL NOISE

We evaluate robustness by randomly corrupting a ratio $\alpha \in \{0.1, 0.2, 0.3, 0.4, 0.5\}$ of CIFAR-20 labels. Final accuracies are reported in Table 3, with training and test dynamics in Figure 12. Training accuracy values closest to $1 - \alpha$ indicate selective fitting of correctly labeled samples.

CE attains nearly perfect training accuracy across all α , showing that it memorizes noisy labels. Consequently, its test accuracy drops sharply (73.6% \rightarrow 40.6% as α increases from 0.1 to 0.5). In contrast, SiMCoding fits mainly the correctly labeled portion, with training accuracy tracking $1 - \alpha$ (e.g., 87.9% at $\alpha = 0.1$ and 49.6% at 0.5). This selective fitting prevents overfitting and yields the strongest test performance, consistently surpassing CE and other baselines. CRATE shows moderate robustness in training but weak generalization (25.4% test at $\alpha = 0.5$). MCR² performs poorly overall, with training accuracy below 30% even at $\alpha = 0.1$ and consistently low test results. Overall, SiMCoding demonstrates the strongest robustness to label noise, limiting fitting to reliable samples and achieving substantially better generalization than CE, CRATE, and MCR².

3.4 ROBUSTNESS AGAINST CATASTROPHIC FORGETTING

[Chan et al. \(2022\)](#) show that optimizing MCR² via iterative gradient ascent naturally induces a multi-layer white-box network, ReduNet. [Wu et al. \(2021\)](#) further adapt ReduNet to incremental learning,

Table 3: Training and test accuracy (%) on CIFAR-20 under different label noise ratios.

Ratio α	Training Accuracy					Test Accuracy				
	0.1	0.2	0.3	0.4	0.5	0.1	0.2	0.3	0.4	0.5
CE	99.74	99.67	99.76	99.79	99.76	73.56	65.94	58.37	47.89	40.57
MCR ²	30.03	23.63	19.71	15.62	13.76	30.37	25.06	22.12	16.37	14.52
CRATE	91.97	90.23	88.94	85.90	84.51	48.00	42.56	35.71	31.21	25.42
SiMCoding	87.86	80.15	71.53	61.25	49.63	77.70	71.26	64.80	57.83	50.70

Table 4: Test accuracy (%) on Task 1 after each training session on MNIST and CIFAR-10.

Algorithm	MNIST					CIFAR-10				
	Task 1	Task 2	Task 3	Task 4	Task 5	Task 1	Task 2	Task 3	Task 4	Task 5
ReduNet	99.95	98.02	95.82	94.94	92.95	78.75	62.35	48.80	44.50	43.07
SiM-ReduNet	99.91	98.05	95.94	94.96	92.25	83.15	63.45	47.97	44.11	40.99

showing reduced catastrophic forgetting. The key difference between MCR² and our SiMCoding lies in subspace dimensionality: MCR² favors maximal dimensions, whereas SiMCoding seeks minimal ones, ideally one-dimensional. This property allows SiMCoding to also yield a ReduNet-like network. We therefore propose SiM-ReduNet, obtained by optimizing SiMCoding through iterative gradient ascent, and show that it offers stronger robustness to catastrophic forgetting and improved performance over ReduNet. A limitation of ReduNet-based incremental learning is memory cost. We adopt a simplified architecture compared to Wu et al. (2021), but still demonstrate that SiMCoding can be directly applied to incremental learning while retaining robustness.

Experimental setup. We evaluate on MNIST and CIFAR-10 under a class-incremental setting, splitting the 10 classes into 5 tasks of 2 classes each. After each task, performance is measured on all classes seen so far. For MNIST, we follow Wu et al. (2021) but reduce the network depth to 50 layers. For CIFAR-10, we omit Gaussian kernel lifting, use a shallower 10-layer network instead of 50, and increase the learning rate from $\eta = 0.5$ to $\eta = 2.5$.

Table 4 reports test accuracy on the first task after each incremental training session. On MNIST, ReduNet and SiM-ReduNet perform similarly, both maintaining high accuracy; SiM-ReduNet is slightly stronger on intermediate tasks but marginally weaker at the final task. On CIFAR-10, forgetting is more pronounced. SiM-ReduNet starts from a higher baseline (e.g., 83.15% vs. 78.75% after Task 1) and holds an advantage in early stages, though ReduNet retains slightly more by the last task. These results confirm that SiMCoding can be effectively extended to incremental learning: SiM-ReduNet closely matches ReduNet’s robustness to catastrophic forgetting and delivers stronger performance on earlier tasks.

4 CONCLUSION

We introduced SiMCoding, a framework that learns representations where each class is characterized by a single, most discriminative component, in contrast to the maximal structural details favoured by MCR². This shift yields strong empirical benefits: SiMCoding consistently matches or surpasses cross-entropy in generalization while showing markedly stronger robustness to label and feature noise. Compared to CRATE and MCR², it demonstrates superior fitting capacity, stability, and inter-class separation across moderate-scale benchmarks. When adapted to incremental learning, SiMCoding also shows robustness against catastrophic forgetting. Despite its computational complexity, these results establish SiMCoding as a simple yet powerful alternative to cross-entropy and related coding-based methods on moderate datasets, with strong potential for robust representation learning and incremental learning.

Empirically, we observed that both CRATE and SiMCoding pursue minimal-dimensional, mutually orthogonal subspaces. Establishing a rigorous connection between the two would be valuable for both theoretical understanding and practical application, which we leave for future work.

REFERENCES

- Kwan Ho Ryan Chan, Yaodong Yu, Chong You, Haozhi Qi, John Wright, and Yi Ma. Redunet: A white-box deep network from the principle of maximizing rate reduction. *Journal of machine learning research*, 23(114):1–103, 2022.
- Hugh Chen, Ian C Covert, Scott M Lundberg, and Su-In Lee. Algorithms to estimate shapley value feature attributions. *Nature Machine Intelligence*, 5(6):590–601, 2023.
- Tianzhe Chu Chu, Shengbang Tong, Tianjiao Ding, Xili Dai, Benjamin Haeffele, René Vidal, and Yi Ma. Image clustering via the principle of rate reduction in the age of pretrained models. International Conference on Learning Representations (ICLR), 2024.
- Thomas M Cover. *Elements of information theory*. John Wiley & Sons, 1999.
- Xili Dai, Shengbang Tong, Mingyang Li, Ziyang Wu, Michael Psenka, Kwan Ho Ryan Chan, Pengyuan Zhai, Yaodong Yu, Xiaojun Yuan, Heung-Yeung Shum, et al. Ctrl: Closed-loop transcription to an ldr via minimaxing rate reduction. *Entropy*, 24(4):456, 2022.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- Matthias Hein and Jean-Yves Audibert. Intrinsic dimensionality estimation of submanifolds in rd. In *Proceedings of the 22nd international conference on Machine learning*, pp. 289–296, 2005.
- Jeremy Howard. Imagenette: A smaller subset of 10 easily classified classes from imagenet, March 2019. URL <https://github.com/fastai/imagenette>.
- Shizhe Hu, Zhengzheng Lou, Xiaoqiang Yan, and Yangdong Ye. A survey on information bottleneck. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- Vignesh Kothapalli. Neural collapse: A review on modelling principles and generalization. *arXiv preprint arXiv:2206.04041*, 2022.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- Yi Ma, Harm Derksen, Wei Hong, and John Wright. Segmentation of multivariate mixed data via lossy data coding and compression. *PAMI*, 2007.
- Jorge Nocedal and Stephen Wright. *Numerical optimization*. Springer Science & Business Media, 2006.
- Druv Pai, Ziyang Wu Wu, Sam Buchanan, Yaodong Yu, and Yi Ma. Masked completion via structured diffusion with white-box transformers. International Conference on Learning Representations, 2023.
- Vardan Pappayan, X. Y. Han, and David L. Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40):24652–24663, sep 2020. doi: 10.1073/pnas.2015509117. URL <https://doi.org/10.1073/pnas.2015509117>.
- Phillip Pope, Chen Zhu, Ahmed Abdelkader, Micah Goldblum, and Tom Goldstein. The intrinsic dimension of images and its impact on learning. *arXiv preprint arXiv:2104.08894*, 2021.
- Jorma Rissanen. *Stochastic complexity in statistical inquiry*, volume 15. World scientific, 1998.
- Stefano Spigler, Mario Geiger, and Matthieu Wyart. Asymptotic learning curves of kernel methods: empirical data vs teacher-student paradigm. *arXiv preprint arXiv:1905.10843*, 2019.

- Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. In *2015 IEEE Information Theory Workshop (ITW)*, pp. 1–5. IEEE, 2015.
- Shengbang Tong, Xili Dai, Yubei Chen, Mingyang Li, Zengyi Li, Brent Yi, Yann LeCun, and Yi Ma. Unsupervised learning of structured representations via closed-loop transcription. *arXiv preprint arXiv:2210.16782*, 2022.
- Shengbang Tong, Xili Dai, Ziyang Wu, Mingyang Li, Brent Yi, and Yi Ma. Incremental learning of structured memory via closed-loop transcription. 2023.
- Peng Wang, Huikang Liu, Druv Pai, Yaodong Yu, Zhihui Zhu, Qing Qu, and Yi Ma. A global geometric analysis of maximal coding rate reduction. *arXiv preprint arXiv:2406.01909*, 2024.
- John Wright and Yi Ma. *High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications*. Cambridge University Press, 2021.
- Ziyang Wu, Christina Baek, Chong You, and Yi Ma. Incremental learning via rate reduction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 1125–1133, 2021.
- Ziyang Wu, Tianjiao Ding, Yifu Lu, Druv Pai, Jingyuan Zhang, Weida Wang, Yaodong Yu, Yi Ma, and Benjamin D Haeffele. Token statistics transformer: Linear-time attention via variational rate reduction. *arXiv preprint arXiv:2412.17810*, 2024.
- Ziyang Wu, Jingyuan Zhang, Druv Pai, XuDong Wang, Chandan Singh, Jianwei Yang, Jianfeng Gao, and Yi Ma. Simplifying dino via coding rate regularization. *arXiv preprint arXiv:2502.10385*, 2025.
- Jinrui Yang, Xianhang Li, Druv Pai, Yuyin Zhou, Yi Ma, Yaodong Yu, and Cihang Xie. Scaling white-box transformers for vision. *Advances in Neural Information Processing Systems*, 37: 36995–37019, 2024.
- Yaodong Yu, Kwan Ho Ryan Chan, Chong You, Chaobing Song, and Yi Ma. Learning diverse and discriminative representations via the principle of maximal coding rate reduction. *Advances in Neural Information Processing Systems*, 33, 2020.
- Yaodong Yu, Sam Buchanan, Druv Pai, Tianzhe Chu, Ziyang Wu, Shengbang Tong, Benjamin Haeffele, and Yi Ma. White-box transformers via sparse rate reduction. *Advances in Neural Information Processing Systems*, 36:9422–9457, 2023.
- Yaodong Yu, Sam Buchanan, Druv Pai, Tianzhe Chu, Ziyang Wu, Shengbang Tong, Hao Bai, Yuexiang Zhai, Benjamin D Haeffele, and Yi Ma. White-box transformers via sparse rate reduction: Compression is all there is? *Journal of Machine Learning Research*, 25(300):1–128, 2024a.
- Yaodong Yu, Tianzhe Chu, Shengbang Tong, Ziyang Wu, Druv Pai, Sam Buchanan, and Yi Ma. Emergence of segmentation with minimalistic white-box transformers. In *Conference on Parsimony and Learning*, volume 234 of *Proceedings of Machine Learning Research*, pp. 72–93. PMLR, 2024b.
- Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision*, pp. 818–833. Springer, 2014.
- Zhihui Zhu, Tianyu Ding, Jinxin Zhou, Xiao Li, Chong You, Jeremias Sulam, and Qing Qu. A geometric analysis of neural collapse with unconstrained features. *Advances in Neural Information Processing Systems*, 34:29820–29834, 2021.

Appendix

LLM usage. ChatGPT was used solely for language polishing of this paper.

A PROOF ABOUT SIMCODING

Notations. We denote by \mathbb{S}_{++}^d the set of $d \times d$ symmetric positive definite matrices, by $\mathbb{R}_{\geq 0}$ the set of non-negative real numbers, and by \mathbb{Z}_{++} the set of positive integers.

A.1 PRELIMINARIES

We begin with a few lemmas that support both the implementation of our algorithm and the subsequent theoretical analysis.

Since $\mathbf{Z}\mathbf{Z}^\top \in \mathbb{R}^{d \times d}$ and $\mathbf{Z}^\top \mathbf{Z} \in \mathbb{R}^{m \times m}$ share identical non-zero eigenvalues, the coding rate admits the following commutative form.

Lemma 1 (Commutative Property (Ma et al., 2007)). *For any $\mathbf{Z} \in \mathbb{R}^{d \times m}$,*

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

When $d < m$, the second expression is computationally more efficient, which is particularly useful for implementing our DARR framework.

In addition, the singular values of \mathbf{Z} (equivalently, the eigenvalues of $\mathbf{Z}\mathbf{Z}^\top$) remain unchanged under orthogonal transformations, yielding the following invariance property.

Lemma 2 (Invariance Property (Ma et al., 2007)). *For any $\mathbf{Z} \in \mathbb{R}^{d \times m}$ and any orthogonal matrices $\mathbf{U} \in \mathbb{R}^{d \times d}$ and $\mathbf{V} \in \mathbb{R}^{m \times m}$,*

$$\mathcal{R}(\mathbf{Z}, \epsilon) = \mathcal{R}(\mathbf{U}\mathbf{Z}\mathbf{V}^\top, \epsilon).$$

Lemma 3 (Upper Bound of MCR² (Yu et al., 2020)). *Let $\mathbf{Z} \in \mathbb{R}^{d \times m}$ denote the complete set of representations for K classes, where the i -th class is represented by $\mathbf{Z}_i \in \mathbb{R}^{d \times m_i}$, i.e., $\mathbf{Z} = \cup_{i=1}^K \mathbf{Z}_i$ and $m = \sum_{i=1}^K m_i$. We have*

$$\Delta \mathcal{R}(\mathbf{Z}, \epsilon) \leq \sum_{i=1}^K \frac{1}{2m} \log \left(\frac{\det^m \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_i \mathbf{Z}_i^\top \right)}{\det^{m_i} \left(\mathbf{I} + \frac{d}{m_i \epsilon^2} \mathbf{Z}_i \mathbf{Z}_i^\top \right)} \right), \quad (3)$$

with equality holds if and only if

$$\mathbf{Z}_i^\top \mathbf{Z}_j = \mathbf{0} \quad \text{for all } 1 \leq i \neq j \leq K.$$

A.2 PROPERTIES OF SIMCODING FRAMEWORK

We now present a more formal version of Theorem 1:

Theorem 2. *Consider the following optimisation objective:*

$$\begin{aligned} \max_{\mathbf{Z}} \Delta \mathcal{R}(\mathbf{Z}, \epsilon) &= \frac{1}{2} \log \det \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}\mathbf{Z}^\top \right) - \sum_{i=1}^K \frac{m_i}{2m} \log \det \left(\mathbf{I} + \frac{d}{m_i \epsilon^2} \mathbf{Z}_i \mathbf{Z}_i^\top \right), \\ \text{s.t. } &\|\mathbf{Z}_i\|_F^2 = m_i \end{aligned} \quad (4)$$

Let $\mathbf{Z}^* = \mathbf{Z}_1^* \cup \dots \cup \mathbf{Z}_K^*$ be the optimal solution to equation 4, and

$$d_i^{\max} = \sqrt{\frac{m_i}{m}} \cdot \frac{d}{\epsilon^2}.$$

Then the following properties hold:

- **Discriminativeness:** Representations from different classes lie in mutually orthogonal low-dimensional linear subspaces:

$$(\mathbf{Z}_i^*)^\top \mathbf{Z}_j^* = 0, \text{ for all } i \neq j.$$

- **Bounded Dimensionality:** Each class-specific subspace satisfies $d_i \leq d_i^{\max}$. The singular values of \mathbf{Z}_i^* fall into one of two patterns:

1. All singular values are equal to $\sqrt{\frac{m_i}{d_i}}$.
2. The largest $d_i - 1$ singular values are equal to σ_H , and the remaining singular value is σ_L , where

$$\sigma_H \in \left(\sqrt{\frac{m_i}{d_i}}, \sqrt{\frac{m_i}{d_i - 1}} \right), \text{ and } \sigma_L > 0.$$

A.3 PROOF OF MAIN RESULTS

We start with presenting two lemmas that will be useful for the proof for Theorem 2.

Lemma 4. Consider the following optimisation problem:

$$\max_m m \cdot f\left(\frac{c}{m}\right) \quad \text{subject to} \quad m \leq \frac{c}{x_p}, \quad (5)$$

where the function $f(x)$ satisfies the following properties:

- (i) $f(0) = 0$,
- (ii) $f'(0) = 0$ and $f'(x) > 0$ for all $x > 0$,
- (iii) There exists $x_p > 0$ such that $f'(x)$ is strictly increasing on $(0, x_p)$ and strictly decreasing on (x_p, ∞) , with $f''(x_p) = 0$.

Let m^* be the optimal solution to equation 5. Then:

$$m^* < \frac{c}{x_p}.$$

Proof of Lemma 4. Let $x = \frac{c}{m}$, and define the auxiliary function

$$g(x) = \frac{f(x)}{x}, \quad \text{so that} \quad g'(x) = \frac{xf'(x) - f(x)}{x^2}.$$

Next, define

$$h(x) = xf'(x) - f(x), \quad \text{which implies} \quad h'(x) = xf''(x).$$

By assumption (iii), $f''(x_p) = 0$, and:

- For $x \in (0, x_p)$, we have $f''(x) > 0$, so $h'(x) > 0$;
- For $x > x_p$, $f''(x) < 0$, so $h'(x) < 0$.

Since $h(0) = 0$, this implies that $h(x)$ increases on $(0, x_p)$, decreases on (x_p, ∞) , and attains its maximum at $x = x_p$. In particular, $h(x_p) \geq 0$, and thus:

$$g'(x_p) = \frac{h(x_p)}{x_p^2} \geq 0.$$

Moreover, since $h'(x) < 0$ for $x > x_p$, it follows that $g'(x)$ is strictly decreasing on $[x_p, \infty)$.

Now suppose $g'(x^*) = 0$ at the optimum. Given that $g'(x_p) \geq 0$ and $g'(x)$ is strictly decreasing for $x \geq x_p$, we must have $x^* > x_p$. Recalling that $x^* = \frac{c}{m^*}$, this yields:

$$\frac{c}{m^*} > x_p \quad \Rightarrow \quad m^* < \frac{c}{x_p}.$$

□

Lemma 5. Given any twice differentiable function $f : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$, an integer $r \in \mathbb{Z}_{++}$, and a constant $c \in \mathbb{R}_+$, consider the optimization problem

$$\begin{aligned} \max_{\mathbf{x}} \quad & \sum_{i=1}^r f(x_i) \\ \text{s.t.} \quad & x_1 \geq x_2 \geq \dots \geq x_r \geq 0, \\ & \sum_{i=1}^r x_i = c, \end{aligned} \tag{6}$$

Let \mathbf{x}^* be an arbitrary global solution to equation 6 and $f(x)$ has the following properties:

- (i) $f'(x)$ satisfies $f'(0) = 0$, is positive for $x > 0$;
- (ii) There exists some $x_p > 0$ such that $f'(x)$ is strictly increasing on $(0, x_p]$ and strictly decreasing on $[x_p, \infty)$.

Let

$$m \in \mathbb{Z}_{++} \quad \text{and} \quad m \leq \min\left\{r, \frac{c}{x_p}\right\},$$

then \mathbf{x}^* must take one of the following two forms:

- $\mathbf{x}^* = \left[\underbrace{\frac{c}{m}, \dots, \frac{c}{m}}_m, \underbrace{0, \dots, 0}_{r-m} \right],$
- $\mathbf{x}^* = \left[\underbrace{x_H, \dots, x_H, x_L}_m, \underbrace{0, \dots, 0}_{r-m} \right], \text{ for some } x_H \in \left(\frac{c}{m}, \frac{c}{m-1} \right) \text{ and } x_L > 0.$

Proof. We consider two cases.

Case 1: $r = 1$.

In this case the problem reduces to

$$\max_{x_1} f(x_1) \quad \text{s.t.} \quad x_1 = c.$$

The unique solution is

$$\mathbf{x}^* = [c],$$

which agrees with the stated result.

Case 2: $r > 1$.

Let

$$\begin{aligned} c_i(x) &= -x_i \leq 0. \\ h(x) &= \sum_{i=1}^r x_i - c = 0. \end{aligned}$$

so the Lagrangian function is

$$\begin{aligned} \mathcal{L}(x_1, \dots, x_r, \alpha_1, \dots, \alpha_r, \beta) &= - \sum_{i=1}^r f(x_i) + \sum_{i=1}^r \alpha_i c_i(x) + \beta h(x) \\ &= - \sum_{i=1}^r f(x_i) - \sum_{i=1}^r \alpha_i x_i + \beta \left(\sum_{i=1}^r x_i - c \right) \end{aligned} \tag{7}$$

Thus, the first-order KKT conditions require that the optimal solution \mathbf{x}^* , α^* , and β^* satisfy:

$$\begin{aligned} f'(x_i^*) &= \alpha_i^* - \beta^*, \quad \forall i \in \{1, \dots, r\}, \\ \alpha_i^* x_i^* &= 0, \quad \forall i \in \{1, \dots, r\}, \\ x_i^* &\geq 0, \quad \forall i \in \{1, \dots, r\}, \\ \alpha_i^* &\geq 0, \quad \forall i \in \{1, \dots, r\}, \\ \sum_{i=1}^r x_i^* &= c. \end{aligned} \tag{8}$$

If $x_i^* > 0$, then $\alpha_i^* = 0$, which implies

$$f'(x_i^*) = -\beta^*.$$

If $x_i^* = 0$, recalling that $f'(0) = 0$, we have

$$f'(x_i^*) = \alpha_i^* - \beta^* = 0.$$

So now we assume there are m entries in \mathbf{x}^* are positive, i.e.,

$$\mathbf{x}^* = [\underbrace{x_1^*, \dots, x_m^*}_m, \underbrace{0, \dots, 0}_{r-m}].$$

and

$$f'(x_i^*) = -\beta^*, \forall i \in \{1, \dots, m\}.$$

Since there exists some $x_p > 0$ such that $f'(x)$ is strictly increasing on $(0, x_p]$ and strictly decreasing on $[x_p, \infty)$ —that is, f' is strictly unimodal—the equation $f'(x) = -\beta^*$ has at most two solutions. Consequently, the optimal vector $\mathbf{x}^* = [x_1^*, \dots, x_m^*]$ can take at most two distinct values. Two cases arise:

1. Uniform Case:

If $x_1^* = x_2^* = \dots = x_m^*$, then

$$mx_i^* = c \implies x_i^* = \frac{c}{m} \quad \text{for all } i,$$

and thus

$$\mathbf{x}^* = [\underbrace{\frac{c}{m}, \dots, \frac{c}{m}}_m, \underbrace{0, \dots, 0}_{r-m}].$$

In this case, it is necessary to note that

$$\frac{c}{m} \geq x_p \iff m \leq \frac{c}{x_p}.$$

We prove it by contradiction. Suppose, on the contrary, that $\frac{c}{m} < x_p$. Then there exists $\sigma > 0$ such that $\frac{c}{m} + \sigma < x_p$. Define

$$\hat{x}_1^* = \frac{c}{m} - \sigma, \quad \hat{x}_2^* = \frac{c}{m} + \sigma.$$

Since $f'(x)$ is strictly increasing on $(0, x_p]$, it follows that

$$f'\left(\frac{c}{m} - \sigma\right) < f'\left(\frac{c}{m} + \sigma\right).$$

This inequality implies the slope between $\frac{c}{m} - \sigma$ and $\frac{c}{m}$ is less than the slope between $\frac{c}{m}$ and $\frac{c}{m} + \sigma$, hence

$$\frac{f\left(\frac{c}{m}\right) - f\left(\frac{c}{m} - \sigma\right)}{\sigma} < \frac{f\left(\frac{c}{m} + \sigma\right) - f\left(\frac{c}{m}\right)}{\sigma}.$$

Multiplying by σ and rearranging terms yields

$$f(\hat{x}_1^*) + f(\hat{x}_2^*) = f\left(\frac{c}{m} - \sigma\right) + f\left(\frac{c}{m} + \sigma\right) > 2f\left(\frac{c}{m}\right) = f(x_1^*) + f(x_2^*),$$

i.e.,

$$f(\hat{x}_1^*) + f(\hat{x}_2^*) + \sum_{i=3}^r f(x_i^*) > \sum_{i=1}^r f(x_i^*).$$

This contradiction shows that the assumed inequality $\frac{c}{m} < x_p$ cannot hold. Hence, we conclude

$$\frac{c}{m} \geq x_p, \quad \text{equivalently,} \quad m \leq \frac{c}{x_p}.$$

We now consider the case where multiple values of m satisfy $m \leq \frac{c}{x_p}$. A natural question is whether the optimum lies on the boundary, i.e., $m^* = \frac{c}{x_p}$. By Lemma 4, it does not: the optimal value strictly satisfies

$$m^* < \frac{c}{x_p}.$$

2. Two-Valued Case:

Otherwise, the set $\{y_1^*, \dots, y_m^*\}$ consists of two distinct values, say x_H and x_L with $x_H > x_L > 0$. Without loss of generality, assume that, after a suitable reordering,

$$x_1^* = \dots = x_{m-1}^* = x_H \quad \text{and} \quad x_m^* = x_L.$$

Then the constraint implies

$$(m-1)x_H + x_L = c.$$

Notice that the average c/m must lie strictly between x_L and x_H ; hence,

$$0 < x_L < \frac{c}{m} \quad \text{and} \quad \frac{c}{m} < x_H < \frac{c}{m-1}.$$

By the unimodality of f' we have

$$x_L < x_p < x_H.$$

Consequently,

$$f''(x_L) > 0 \quad \text{and} \quad f''(x_H) < 0.$$

In this case, $x_H > x_p$. To make it, we may push a stronger condition, i.e.,

$$x_H > \frac{c}{m} \geq x_p \Rightarrow m \leq \frac{c}{x_p}.$$

In this case, if we fix x_L , and only consider the rest $(m-1)x_H$, according to Lemma 4, the optimal value satisfies

$$m^* - 1 < \frac{c - x_L}{x_p} \Rightarrow m^* < \frac{c}{x_p} + \frac{x_p - x_L}{x_p}.$$

Combining the constraint $m \leq \frac{c}{x_p}$ with the optimal $m^* < \frac{c}{x_p} + \frac{x_p - x_L}{x_p}$, we conclude that

$$m^* \leq \frac{c}{x_p}.$$

Second-Order Analysis:

To ensure that the candidate \mathbf{y}^* is indeed optimal, the second-order necessary conditions must hold (see, e.g., (Nocedal & Wright, 2006, Theorem 12.5)). For any feasible direction $\mathbf{v} = (v_1, \dots, v_m)$ satisfying

$$\sum_{i=1}^m v_i = 0,$$

the Hessian of the Lagrangian must be negative semidefinite on the tangent space, i.e.,

$$\sum_{i=1}^m f''(x_i^*) v_i^2 \leq 0.$$

Suppose, for the sake of contradiction, that more than one index takes the lower value x_L . For example, if $x_{m-1}^* = x_m^* = x_L$, consider the perturbation direction defined by

$$v_{m-1} = 1, \quad v_m = -1, \quad \text{and } v_i = 0 \text{ for } i \neq m-1, m.$$

Clearly, $\sum_{i=1}^m v_i = 1 + (-1) = 0$. Then the second-order condition yields

$$\sum_{i=1}^m f''(x_i^*) v_i^2 = f''(x_{m-1}^*) + f''(x_m^*) = 2f''(x_L) \leq 0.$$

However, since $x_L < x_p$ we have $f''(x_L) > 0$, so that $2f''(x_L) > 0$, which contradicts the second-order necessary condition. Hence, it is necessary that exactly one index takes the lower value x_L , with all remaining $m-1$ indices equal to x_H .

Then the solution in this case is

$$\mathbf{x}^* = \left[\underbrace{x_H, \dots, x_H}_m, \underbrace{x_L, 0, \dots, 0}_{r-m} \right],$$

where, from the relation $(m-1)x_H + x_L = c$, we obtain

$$(m-1)x_H < c \quad \text{and} \quad mx_H > c.$$

Thus,

$$x_H \in \left(\frac{c}{m}, \frac{c}{m-1} \right) \quad \text{and} \quad x_L > 0.$$

Step 5: Conclusion.

Combining the results from the two cases, we deduce that given

$$m \leq \min \left\{ r, \frac{c}{x_p} \right\} \leq \frac{c}{x_p},$$

any global solution \mathbf{x}^* to problem equation 6 must have the form

$$\mathbf{x}^* = \left[\underbrace{x_1^*, \dots, x_m^*}_m, \underbrace{0, \dots, 0}_{r-m} \right],$$

where either

$$x_1^* = \dots = x_m^* = \frac{c}{m},$$

or

$$\mathbf{x}^* = \left[\underbrace{x_H, \dots, x_H}_m, \underbrace{x_L, 0, \dots, 0}_{r-m} \right],$$

with

$$x_H \in \left(\frac{c}{m}, \frac{c}{m-1} \right) \quad \text{and} \quad x_L > 0.$$

This completes the proof. \square

Proof of Theorem 2. This proof follows a structure similar to that of Theorem A.6 in [Yu et al. \(2020\)](#).

Specifically, without loss of generality, let

$$\mathbf{Z}^* = [\mathbf{Z}_1^*, \dots, \mathbf{Z}_K^*]$$

be an optimal solution to problem equation 2. To prove that the matrices \mathbf{Z}_i^* , for $i \in \{1, \dots, K\}$, are pairwise orthogonal, assume for the sake of contradiction that

$$(\mathbf{Z}_i^*)^\top \mathbf{Z}_j^* \neq \mathbf{0}$$

for some indices $1 \leq i < j \leq K$.

By applying Lemma 3, the strict inequality in equation 3 holds for the optimal solution \mathbf{Z}^* . That is,

$$\Delta \mathcal{L}(\mathbf{Z}^*, \alpha) < \sum_{i=1}^K \frac{1}{2m} \log \left(\frac{\det^m \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_i^* (\mathbf{Z}_i^*)^\top \right)}{\det^{m_i} \left(\mathbf{I} + \frac{d}{m_i \epsilon^2} \mathbf{Z}_i^* (\mathbf{Z}_i^*)^\top \right)} \right). \quad (9)$$

On the other hand, let $\text{rank}(\mathbf{Z}_i) = d_i$, there exist matrices

$$\{\mathbf{U}'_i \in \mathbb{R}^{d \times d_i}\}_{i=1}^K$$

such that the columns of the concatenated matrix $[\mathbf{U}'_1, \dots, \mathbf{U}'_K]$ form an orthonormal set (i.e., for any $i \neq j$, we have $\mathbf{U}'_i (\mathbf{U}'_j)^\top = \mathbf{0}$, and each column has unit l_2 -norm).

Let the compact singular value decomposition (SVD) of \mathbf{Z}_i^* be given by

$$\mathbf{Z}_i^* = \mathbf{U}_i^* \Sigma_i^* (\mathbf{V}_i^*)^\top,$$

and define

$$\mathbf{Z}' = [\mathbf{Z}'_1, \dots, \mathbf{Z}'_K], \quad \text{where} \quad \mathbf{Z}'_i = \mathbf{U}'_i \Sigma_i^* (\mathbf{V}_i^*)^\top.$$

Then, for all $1 \leq i < j \leq K$, we have

$$\begin{aligned} (\mathbf{Z}'_i)^\top \mathbf{Z}'_j &= \mathbf{V}_i^* \Sigma_i^* (\mathbf{U}'_i)^\top \mathbf{U}'_j \Sigma_j^* (\mathbf{V}_j^*)^\top \\ &= \mathbf{V}_i^* \Sigma_i^* \mathbf{0} \Sigma_j^* (\mathbf{V}_j^*)^\top \\ &= \mathbf{0}. \end{aligned}$$

That is, the matrices $\mathbf{Z}'_1, \dots, \mathbf{Z}'_K$ are pairwise orthogonal.

Applying Lemma 3 to \mathbf{Z}' yields

$$\begin{aligned} \Delta \mathcal{R}(\mathbf{Z}', \alpha) &= \sum_{i=1}^K \frac{1}{2m} \log \left(\frac{\det^m \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}'_i (\mathbf{Z}'_i)^\top \right)}{\det^{m_i} \left(\mathbf{I} + \frac{d}{m_i \epsilon^2} \mathbf{Z}'_i (\mathbf{Z}'_i)^\top \right)} \right) \\ &= \sum_{i=1}^K \frac{1}{2m} \log \left(\frac{\det^m \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_i^* (\mathbf{Z}_i^*)^\top \right)}{\det^{m_i} \left(\mathbf{I} + \frac{d}{m_i \epsilon^2} \mathbf{Z}_i^* (\mathbf{Z}_i^*)^\top \right)} \right), \end{aligned} \quad (10)$$

where the second equality follows from Lemma 2.

Comparing equation 9 and equation 10 shows that

$$\Delta \mathcal{R}(\mathbf{Z}', \alpha) > \Delta \mathcal{L}(\mathbf{Z}^*, \alpha),$$

which contradicts the optimality of \mathbf{Z}^* . Therefore, we must have

$$(\mathbf{Z}_i^*)^\top \mathbf{Z}_j^* = \mathbf{0} \quad \text{for all } 1 \leq i < j \leq K.$$

Therefore, by Lemma 2 we have

$$\Delta \mathcal{R}(\mathbf{Z}^*, \alpha) = \sum_{i=1}^K \frac{1}{2m} \log \left(\frac{\det^m \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_i^* (\mathbf{Z}_i^*)^\top \right)}{\det^{m_i} \left(\mathbf{I} + \frac{d}{m_i \epsilon^2} \mathbf{Z}_i^* (\mathbf{Z}_i^*)^\top \right)} \right). \quad (11)$$

We now prove the result concerning the singular values of \mathbf{Z}_i^* . To begin with, we claim that the following result holds:

$$\begin{aligned} \mathbf{Z}_i^* \in \arg \max_{\mathbf{Z}_i} \log \left(\frac{\det^m \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_i \mathbf{Z}_i^\top \right)}{\det^{m_i} \left(\mathbf{I} + \frac{d}{m_i\epsilon^2} \mathbf{Z}_i \mathbf{Z}_i^\top \right)} \right), \\ \text{s.t. } \|\mathbf{Z}_i\|_F^2 = m_i. \end{aligned} \quad (12)$$

To see why equation 12 holds, suppose that there exists $\tilde{\mathbf{Z}}_i$ satisfying $\|\tilde{\mathbf{Z}}_i\|_F^2 = m_i$ and $\text{rank}(\tilde{\mathbf{Z}}_i) \leq d_i$, such that

$$\log \left(\frac{\det^m \left(\mathbf{I} + \frac{d}{m\epsilon^2} \tilde{\mathbf{Z}}_i \tilde{\mathbf{Z}}_i^\top \right)}{\det^{m_i} \left(\mathbf{I} + \frac{d}{m_i\epsilon^2} \tilde{\mathbf{Z}}_i \tilde{\mathbf{Z}}_i^\top \right)} \right) > \log \left(\frac{\det^m \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_i^* (\mathbf{Z}_i^*)^\top \right)}{\det^{m_i} \left(\mathbf{I} + \frac{d}{m_i\epsilon^2} \mathbf{Z}_i^* (\mathbf{Z}_i^*)^\top \right)} \right). \quad (13)$$

Denote the compact SVD of $\tilde{\mathbf{Z}}_i$ by

$$\tilde{\mathbf{Z}}_i = \tilde{\mathbf{U}}_i \tilde{\Sigma}_i \tilde{\mathbf{V}}_i^\top,$$

and define

$$\mathbf{Z}' = [\mathbf{Z}_1^*, \dots, \mathbf{Z}_{i-1}^*, \mathbf{Z}_i', \mathbf{Z}_{i+1}^*, \dots, \mathbf{Z}_K^*], \quad \text{where } \mathbf{Z}_i' := \tilde{\mathbf{U}}_i \tilde{\Sigma}_i \tilde{\mathbf{V}}_i^\top.$$

Note that $\|\mathbf{Z}_i'\|_F^2 = m_i$, $\text{rank}(\mathbf{Z}_i') \leq d_i$, and $(\mathbf{Z}_i')^\top \mathbf{Z}_j' = \mathbf{0}$ for all $i \neq j$. Hence, \mathbf{Z}' is a feasible solution to problem equation 4 with its components being pairwise orthogonal.

By invoking Lemma 3, Lemma 2, and using inequality equation 13, we have

$$\begin{aligned} \Delta \mathcal{R}(\mathbf{Z}', \alpha) &= \frac{1}{2m} \log \left(\frac{\det^m \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_i' (\mathbf{Z}_i')^\top \right)}{\det^{m_i} \left(\mathbf{I} + \frac{d}{m_i\epsilon^2} \mathbf{Z}_i' (\mathbf{Z}_i')^\top \right)} \right) + \sum_{j \neq i} \frac{1}{2m} \log \left(\frac{\det^m \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)}{\det^{m_j} \left(\mathbf{I} + \frac{d}{m_j\epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)} \right) \\ &= \frac{1}{2m} \log \left(\frac{\det^m \left(\mathbf{I} + \frac{d}{m\epsilon^2} \tilde{\mathbf{Z}}_i \tilde{\mathbf{Z}}_i^\top \right)}{\det^{m_i} \left(\mathbf{I} + \frac{d}{m_i\epsilon^2} \tilde{\mathbf{Z}}_i \tilde{\mathbf{Z}}_i^\top \right)} \right) + \sum_{j \neq i} \frac{1}{2m} \log \left(\frac{\det^m \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)}{\det^{m_j} \left(\mathbf{I} + \frac{d}{m_j\epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)} \right) \\ &> \frac{1}{2m} \log \left(\frac{\det^m \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_i^* (\mathbf{Z}_i^*)^\top \right)}{\det^{m_i} \left(\mathbf{I} + \frac{d}{m_i\epsilon^2} \mathbf{Z}_i^* (\mathbf{Z}_i^*)^\top \right)} \right) + \sum_{j \neq i} \frac{1}{2m} \log \left(\frac{\det^m \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)}{\det^{m_j} \left(\mathbf{I} + \frac{d}{m_j\epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)} \right) \\ &= \sum_{i=1}^K \frac{1}{2m} \log \left(\frac{\det^m \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_i^* (\mathbf{Z}_i^*)^\top \right)}{\det^{m_i} \left(\mathbf{I} + \frac{d}{m_i\epsilon^2} \mathbf{Z}_i^* (\mathbf{Z}_i^*)^\top \right)} \right). \end{aligned}$$

Combining it with equation 11 shows $\Delta \mathcal{R}(\mathbf{Z}', \alpha) > \Delta \mathcal{R}(\mathbf{Z}^*, \alpha)$, contradicting the optimality of \mathbf{Z}^* . Therefore, the result in equation 12 holds.

Observe that the optimization problem in equation 12 depends on \mathbf{Z}_i only through its singular values. More precisely, let

$$\sigma_i := [\sigma_{1,i}, \dots, \sigma_{\min\{m_i, d\}, i}]$$

denote the singular values of \mathbf{Z}_i and

$$\det \left(\mathbf{I} + \gamma \mathbf{Z} \mathbf{Z}^\top \right) = \prod_{p=1}^n \left(1 + \gamma \sigma_p^2 \right).$$

Then, we have

$$\log \left(\frac{\det^m \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_i \mathbf{Z}_i^\top \right)}{\det^{m_i} \left(\mathbf{I} + \frac{d}{m_i\epsilon^2} \mathbf{Z}_i \mathbf{Z}_i^\top \right)} \right) = \sum_{p=1}^{\min\{m_i, d\}} \log \left(\frac{\left(1 + \frac{d}{m\epsilon^2} \sigma_{p,i}^2 \right)^m}{\left(1 + \frac{d}{m_i\epsilon^2} \sigma_{p,i}^2 \right)^{m_i}} \right).$$

Moreover, note that

$$\|\mathbf{Z}_i\|_F^2 = \sum_{p=1}^{\min\{m_i, d\}} \sigma_{p,i}^2.$$

Using these relations, the optimization problem in equation 12 is equivalent to

$$\begin{aligned} \max_{\sigma_i \in \mathcal{R}_{\geq 0}^{\min\{m_i, d\}}} \quad & \sum_{p=1}^{\min\{m_i, d\}} \log \left(\frac{\left(1 + \frac{d}{m\epsilon^2} \sigma_{p,i}^2\right)^m}{\left(1 + \frac{d}{m_i\epsilon^2} \sigma_{p,i}^2\right)^{m_i}} \right) \\ \text{s.t.} \quad & \sum_{p=1}^{\min\{m_i, d\}} \sigma_{p,i}^2 = m_i. \end{aligned} \quad (14)$$

Let

$$r_i = \min\{m_i, d\} \quad \text{and} \quad \sigma_i^* = [\sigma_{1,i}^*, \dots, \sigma_{r_i,i}^*]$$

be an optimal solution to equation 14. Without loss of generality, assume that the entries of σ_i^* are sorted in descending order.

Hence,

$$\begin{aligned} [\sigma_{1,i}^*, \dots, \sigma_{r_i,i}^*] = \quad & \arg \max_{[\sigma_{1,i}, \dots, \sigma_{r_i,i}] \in \mathcal{R}_{\geq 0}^{r_i}} \sum_{p=1}^{r_i} \log \left(\frac{\left(1 + \frac{d}{m\epsilon^2} \sigma_{p,i}^2\right)^m}{\left(1 + \frac{d}{m_i\epsilon^2} \sigma_{p,i}^2\right)^{m_i}} \right) \\ \text{s.t.} \quad & \sum_{p=1}^{r_i} \sigma_{p,i}^2 = m_i, \quad \sigma_{1,i} \geq \dots \geq \sigma_{r_i,i} \geq 0. \end{aligned} \quad (15)$$

To solve the problem equation 15 we define a new function as

$$\begin{aligned} f(x; d, \alpha, m_i, m) &= \log \left(\frac{\left(1 + \frac{d}{m\epsilon^2} x\right)^m}{\left(1 + \frac{d}{m_i\epsilon^2} x\right)^{m_i}} \right) \\ &= m \log \left(1 + \frac{d}{m\epsilon^2} x\right) - m_i \log \left(1 + \frac{d}{m_i\epsilon^2} x\right) \\ \text{s.t.} \quad & x \geq 0 \end{aligned}$$

We compute the first derivative of f with respect to x , which is given by

$$\begin{aligned} f'(x) &= \frac{dm}{m\epsilon^2 + dx} - \frac{dm_i}{m_i\epsilon^2 + dx} \\ &= \frac{d^2(m - m_i)x}{(m\epsilon^2 + dx)(m_i\epsilon^2 + dx)}. \end{aligned}$$

Clearly $f'(0) = 0$.

We compute the second derivative of f with respect to x , which is given by

$$f''(x) = \frac{d^2(m - m_i)(mm_i\epsilon^4 - d^2x^2)}{(m\epsilon^2 + dx)^2(m_i\epsilon^2 + dx)^2}.$$

Let

$$\mathbf{x}_p = \frac{\sqrt{mm_i}\epsilon^2}{d}, \quad f''(\mathbf{x}_p) = 0.$$

The optimisation problem equation 15 satisfies

- (i) $f'(x)$ satisfies $f'(0) = 0$, and is positive for $x > 0$;

- (ii) There exists some $x_p > 0$ such that $f'(x)$ is strictly increasing on $(0, x_p]$ and strictly decreasing on $[x_p, +\infty)$;

Therefore, we may apply Lemma 5 and conclude that let

$$d_i \leq \min \left\{ r_i, \frac{m_i}{x_p} \right\} = \min \left\{ r_i, \sqrt{\frac{m_i}{m}} \frac{d}{\epsilon^2} \right\}$$

Hence for simplicity, we have

$$d_i \leq \sqrt{\frac{m_i}{m}} \frac{d}{\epsilon^2}.$$

The unique optimal solution to equation 15 is one of the following two forms:

- $\mathbf{x}^* = \left[\underbrace{\frac{m_i}{d_i}, \dots, \frac{m_i}{d_i}}_{d_i}, \underbrace{0, \dots, 0}_{r_i - d_i} \right],$
- $\mathbf{x}^* = \left[\underbrace{x_H, \dots, x_H, x_L}_{d_i}, \underbrace{0, \dots, 0}_{r_i - d_i} \right],$ for some $x_H \in \left(\frac{m_i}{d_i}, \frac{m_i}{d_i - 1} \right)$ and $x_L > 0$.

In fact, if the optimal solution is the first form, according to Lemma 5, $d_i < \sqrt{\frac{m_i}{m}} \frac{d}{\epsilon^2}$.

Equivalently, we have either

- $[\sigma_{1,i}^*, \dots, \sigma_{r_i,i}^*] = \left[\underbrace{\sqrt{\frac{m_i}{d_i}}, \dots, \sqrt{\frac{m_i}{d_i}}}_{d_i}, \underbrace{0, \dots, 0}_{r_i - d_i} \right],$ or
- $[\sigma_{1,i}^*, \dots, \sigma_{r_i,i}^*] = \left[\underbrace{\sigma_H, \dots, \sigma_H, \sigma_L}_{d_i}, \underbrace{0, \dots, 0}_{r_i - d_i} \right],$ for some $\sigma_H \in \left(\sqrt{\frac{m_i}{d_i}}, \sqrt{\frac{m_i}{d_i - 1}} \right)$ and $\sigma_L > 0$,

as claimed. □

B ADDITIONAL EXPERIMENTS

B.1 COMPLEMENTARY IMPLEMENTATION DETAILS

Architecture. For the white-box transformer architecture CRATE, we adapt the model complexity to each dataset:

- **MNIST:** A compact grayscale model with a small patch size to mitigate overfitting.
- **CIFAR family:** A wider model tailored to 32×32 color images.
- **ImageNette:** A deeper model suited for higher-resolution 128×128 images.

```
def CRATE_tiny_mnist(num_classes=10):
    # MNIST: 28x28 grayscale
    return CRATE(
        image_size=28,
        patch_size=7,          # 4x4 = 16 tokens
        num_classes=num_classes,
        dim=64,                # model width
        depth=6,               # number of layers
        heads=4,               # attention heads
        dim_head=16,           # 64 / 4
        dropout=0.0,
        emb_dropout=0.0,
        channels=1             # grayscale input
    )

def CRATE_tiny_cifar(num_classes=10):
    # CIFAR: 32x32 RGB
    return CRATE(
        image_size=32,
        patch_size=8,          # 4x4 = 16 tokens
        num_classes=num_classes,
        dim=512,
        depth=8,
        heads=8,
        dim_head=64,           # 512 / 8
        dropout=0.0,
        emb_dropout=0.0,
        channels=3
    )

def CRATE_tiny_imagenette(num_classes=10):
    # ImageNette: 128x128 RGB
    return CRATE(
        image_size=128,
        patch_size=8,          # 16x16 = 64 tokens
        num_classes=num_classes,
        dim=512,
        depth=8,
        heads=8,
        dim_head=64,           # 512 / 8
        dropout=0.0,
        emb_dropout=0.0,
        channels=3
    )
```

Data Augmentation. For the MNIST dataset, we do not apply any data augmentation. For CIFAR family, we follow prior works (Yu et al., 2020; Chan et al., 2022) and adopt the following augmentation pipeline:

```

import torchvision.transforms as transforms
TRANSFORM = transforms.Compose([
    transforms.RandomResizedCrop(32, padding=8),
    transforms.RandomHorizontalFlip(),
    transforms.ToTensor()])

```

For CIFAR-100, we construct the final training set by duplicating the original data five times, which empirically improves SiMCoding fitting. We conjecture that this arises from the limited number of samples per class (only 500) in the original dataset, which may prevent SiMCoding from fully capturing class structure. A more rigorous explanation remains open for future work.

For ImageNette, we use the following transformations for training and test:

```

# Training
train_tf = transforms.Compose([
    transforms.RandomResizedCrop(128),
    transforms.RandomHorizontalFlip(),
    transforms.ToTensor(),
    AddGaussianNoise(mean=0.0, std=args.gauss_std, clip=True)
    if args.gauss_std > 0 else transforms.Lambda(lambda x: x),
    transforms.Normalize(mean=[0.485, 0.456, 0.406],
        std=[0.229, 0.224, 0.225]),
])

# Test
val_tf = transforms.Compose([
    transforms.Resize(160),
    transforms.CenterCrop(128),
    transforms.ToTensor(),
    transforms.Normalize(mean=[0.485, 0.456, 0.406],
        std=[0.229, 0.224, 0.225]),
])

```

Training Setting. For CE, MCR², and SiMCoding, we adopt the following training settings across datasets. For MNIST, CIFAR-10, and CIFAR-20, models are trained using stochastic gradient descent (SGD) with an initial learning rate of 0.01, momentum of 0.9, and weight decay of 5×10^{-4} . Training is performed for 500 epochs with a batch size of 1000, and the learning rate is decayed by a factor of 10 every 200 epochs. For CIFAR-100, we found that a learning rate of 0.01 was too small for convergence, so we increased it to 0.1 while keeping the same decay schedule. For ImageNette, due to the larger input resolution (128×128), we reduce the batch size to 400 while keeping the initial learning rate at 0.01.

For CRATE on all datasets, we adapt the training strategy described in [Yu et al. \(2024a\)](#). The batch size is kept the same as in the other methods for each dataset to ensure a fair comparison. Training is performed using the AdamW optimizer with an initial learning rate of $\eta_0 = 4 \times 10^{-4}$, momentum parameters $\beta = (0.9, 0.999)$, and a weight decay of 0.05. We adopt the label smoothing cross-entropy loss with a smoothing parameter of $\alpha = 0.1$. The models are trained for a total of 500 epochs.

We use a per-iteration LambdaLR scheduler with a linear warmup of 100 steps followed by cosine annealing. For global optimizer step t (starting at 0) and total epochs T , the learning-rate multiplier $\lambda(t)$ is defined as

$$\lambda(t) = \min\left(\frac{t+1}{100+\epsilon}, \frac{1}{2}(1 + \cos(\pi t / \max(1, T)))\right),$$

so that the actual learning rate is $\eta(t) = \eta_0 \lambda(t)$. The scheduler is updated at every training iteration.

Note that in Figure 10, the features learned by MCR² do not fully span the representation space, as the classes are not separated at all. This underscores the limitation of MCR²'s strategy of maximizing structural details.

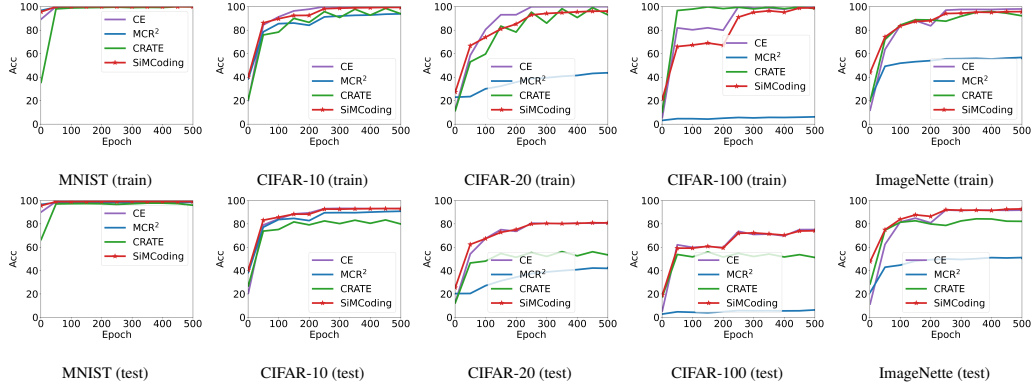


Figure 5: Training (top row) and test (bottom row) accuracy versus epoch on five datasets.

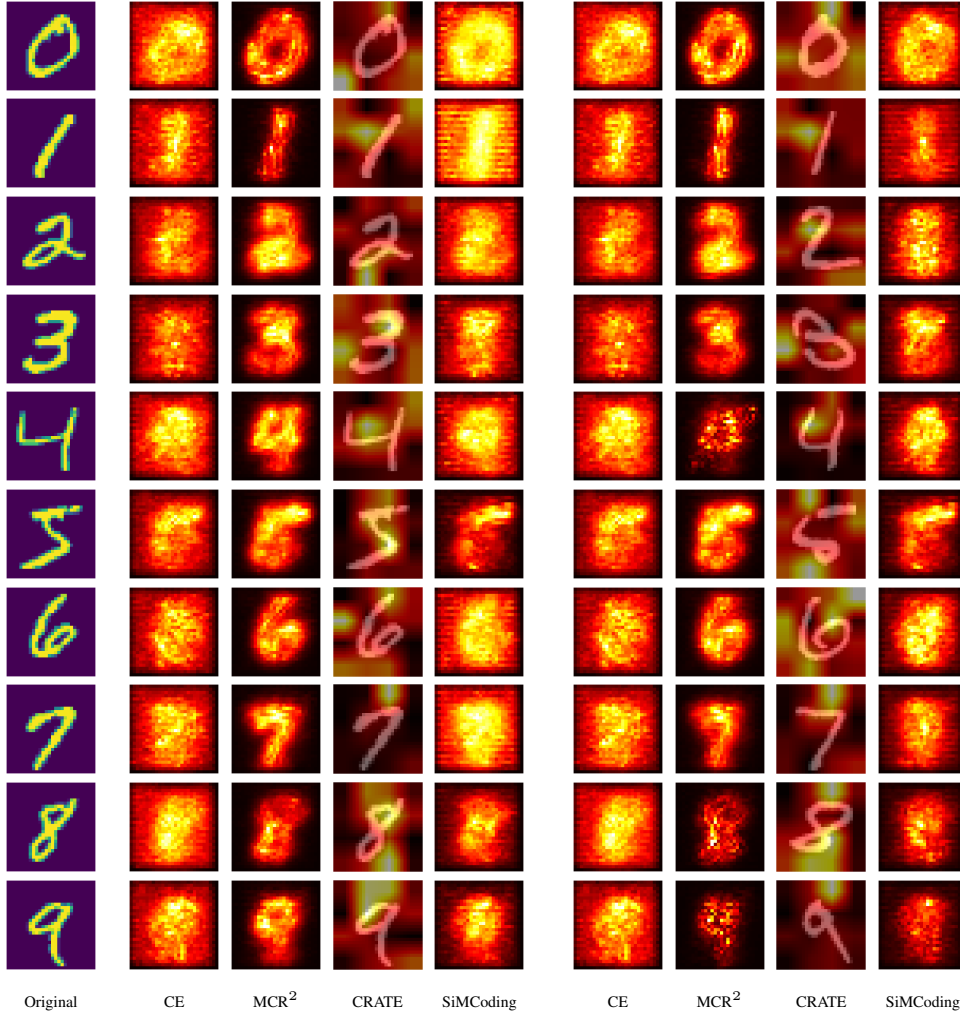


Figure 6: Saliency maps for digits 0 to 9. The first column displays a representative original image for each class from the MNIST training set. The two main blocks show the mean saliency maps for each method (CE, MCR², CRATE and SiMCoding) on the training set (middle) and the test set (right).

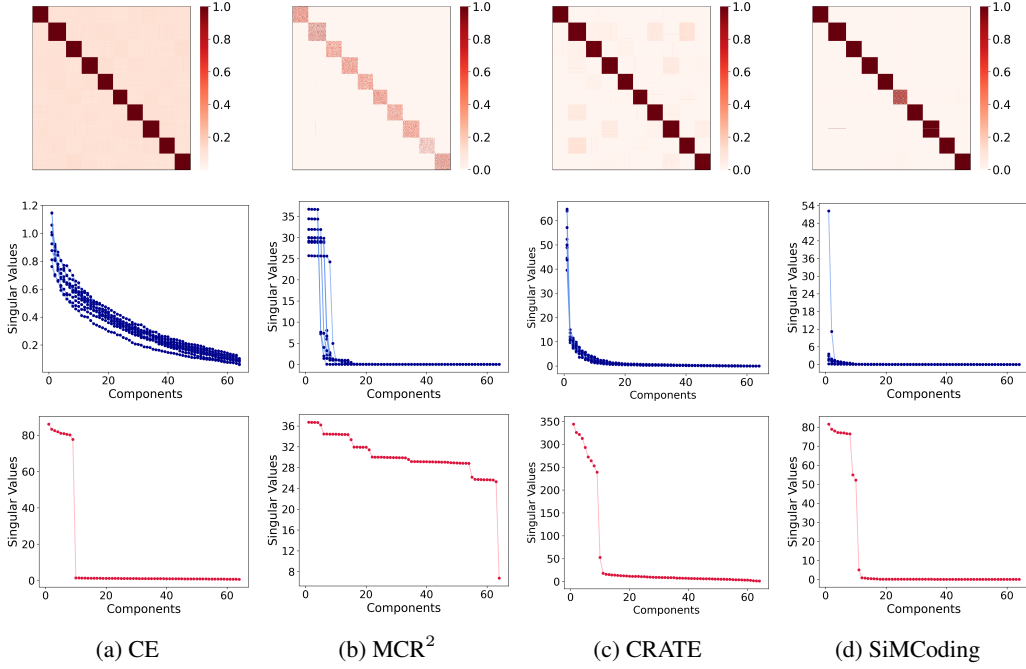


Figure 7: Feature structure analysis on MNIST under different training objectives. Top row: heatmaps of $|\mathbf{Z}\mathbf{Z}^\top|$ with samples sorted by class. Middle row: per-class singular value spectra of \mathbf{Z}_i . Bottom row: overall singular value spectrum of \mathbf{Z} (with method labels).

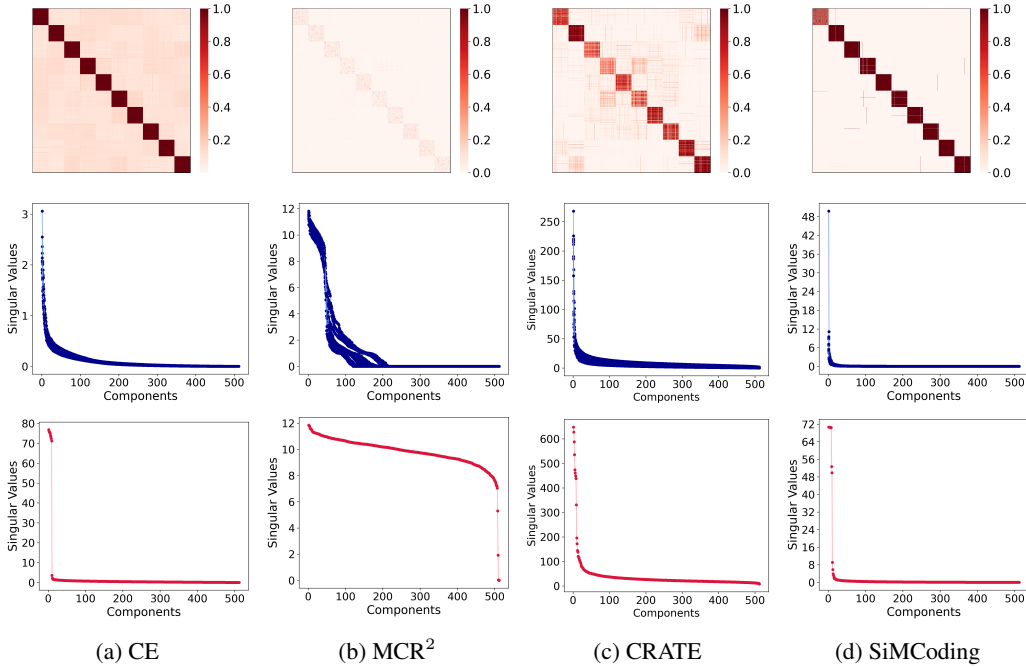


Figure 8: Feature structure analysis on CIFAR10 under different training objectives. Top row: heatmaps of $|\mathbf{Z}\mathbf{Z}^\top|$ with samples sorted by class. Middle row: per-class singular value spectra of \mathbf{Z}_i . Bottom row: overall singular value spectrum of \mathbf{Z} .

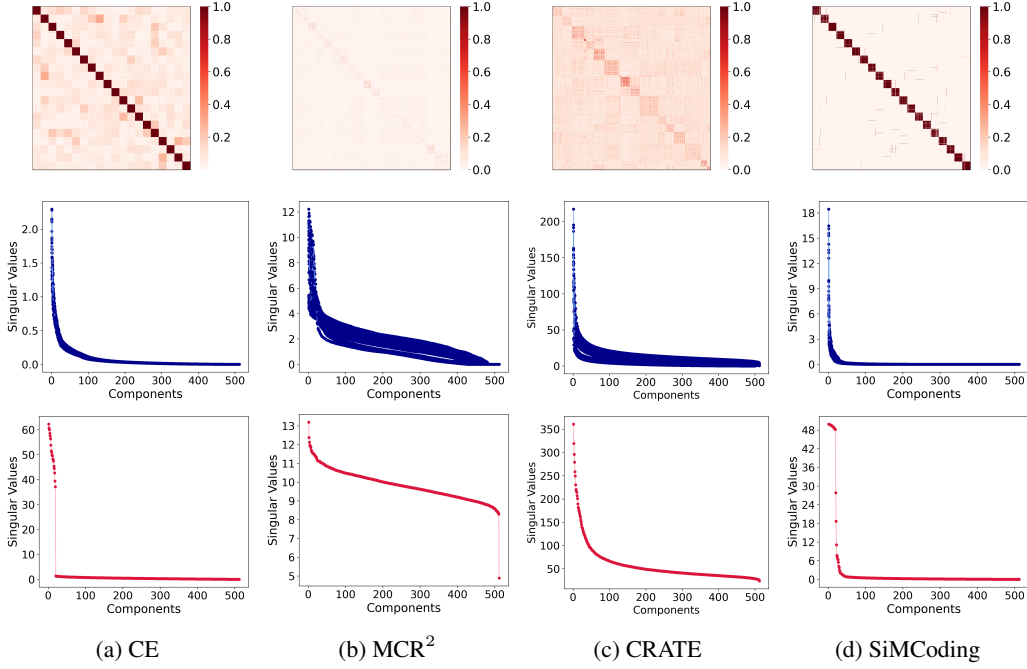


Figure 9: Feature structure analysis on CIFAR20 under different training objectives. Top row: heatmaps of $|\mathbf{Z}\mathbf{Z}^\top|$ with samples sorted by class. Middle row: per-class singular value spectra of \mathbf{Z}_i . Bottom row: overall singular value spectrum of \mathbf{Z} .

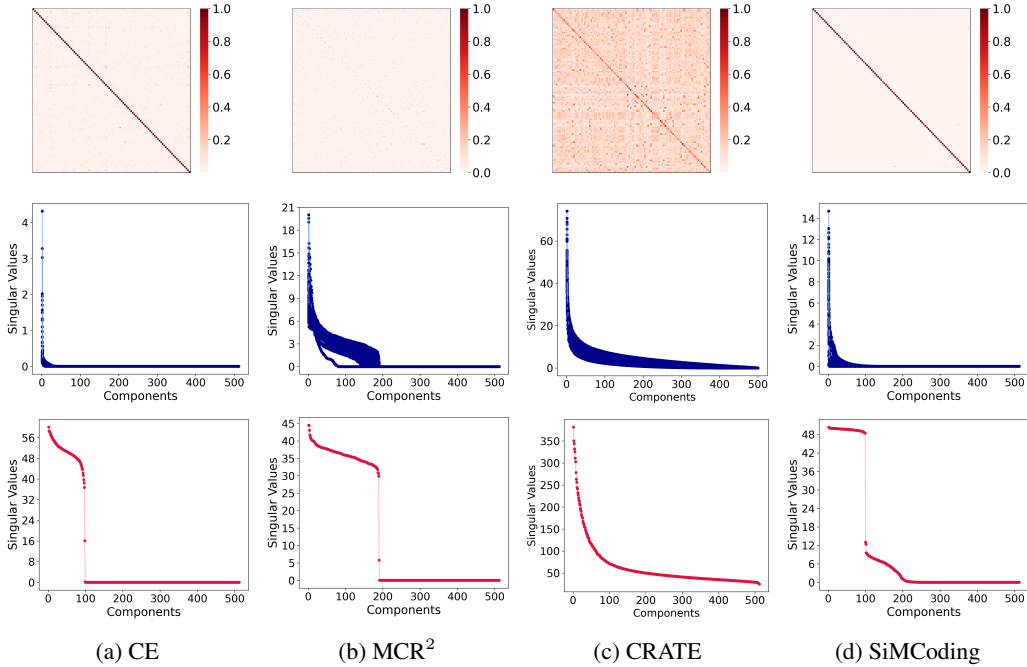


Figure 10: Feature structure analysis on CIFAR100 under different training objectives. Top row: heatmaps of $|\mathbf{Z}\mathbf{Z}^\top|$ with samples sorted by class. Middle row: per-class singular value spectra of \mathbf{Z}_i . Bottom row: overall singular value spectrum of \mathbf{Z} .

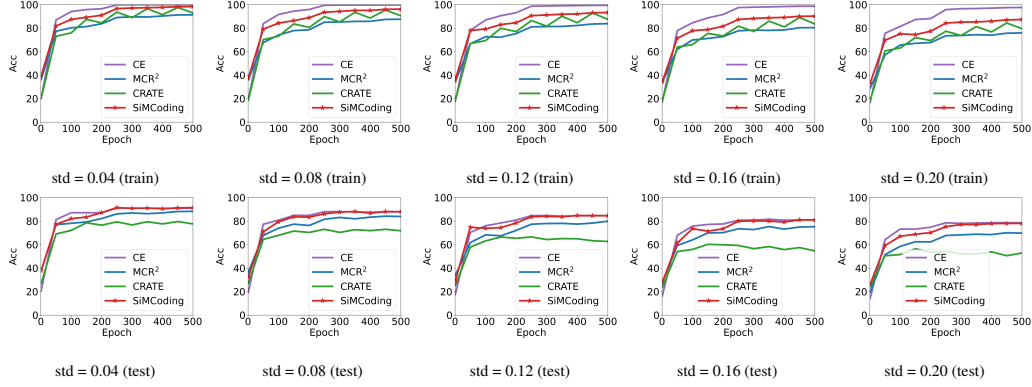


Figure 11: Training (top row) and test (bottom row) accuracy versus epoch on CIFAR-10 with different levels of feature noise. Columns correspond to increasing noise standard deviation (std).

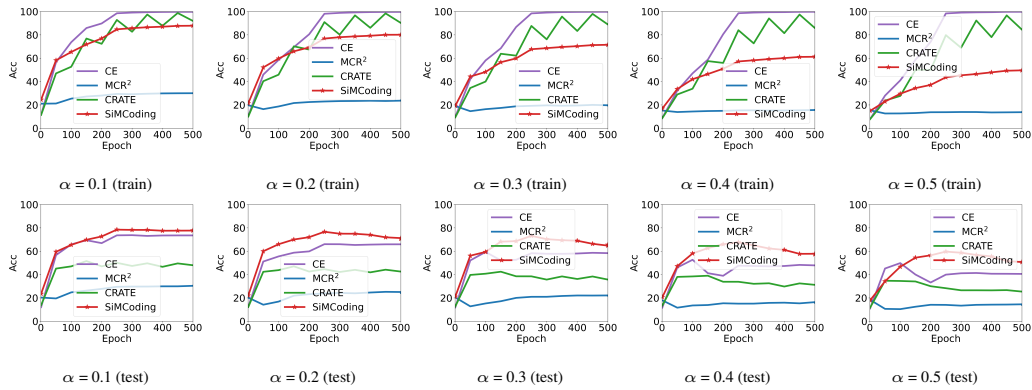


Figure 12: Training (top row) and test (bottom row) accuracy versus epoch on CIFAR-20 with different levels of label noise. Columns correspond to increasing label noise ratios.