How Mask Matters: Towards Theoretical Understandings of Masked Autoencoders

Anonymous Author(s) Affiliation Address email

Abstract

14 1 Introduction

Recently, self-supervised learning has been proposed as an effective framework to learn meaningful 15 representations without access to labels. There are lots of works raising different surrogate objectives 16 to improve the performance of the learned representation on downstream tasks. Among them, masked 17 autoencoding defines self-supervised signals by masking some patches of samples and learns a 18 surrogate loss based on the reconstruction task of the missing message. Bert [6] is a popular self-19 supervised framework based on masked autoencoding in natural language process and achieves 20 impressive performance. Recently, a paradigm of masked autoencoding named Masked Autoencoder 21 (MAE) [11] obtains empirical breakthroughs on various computer vision tasks and achieves state-of-22 the-art performance. However, despite its promising empirical success, a theoretical understanding of 23 why MAE can learn representations that can improve the performance of downstream tasks is still 24 limited. 25

The basic idea of MAE is quite simple, which masks several patches of an image and encodes the unmasked patches with the Vision Transformer (ViT) [7], then tries to reconstruct the missing patches with a decoder via a training objective of mean square error. Intuitively, the reconstruction task is an instance-level task instead of a class-level task. However, the empirical success of MAE claims that the reconstruction target can learn helpful representations for downstream classification tasks. So how to understand the training process of MAE can bridge the samples in the same class is quite important.

[2] is the first theoretical work to understand MAE, which explains MAE from a perspective of
 an integral kernel. However, there still exists a vacancy in their analysis. For example, their work
 doesn't analyze the core component of MAE, *i.e.*, the masking technique. In this paper, we try to
 further understand the training process of MAE and its components. We prove that MAE loss is

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.



Figure 1: For each input x, Masked Autoencoder (MAE) first draws a random patchwise mask m, and obtain two complementary views x_1, x_2 . Afterward, MAE utilizes an encoder-decoder architecture to reconstruct x_2 from x_1 .

³⁷ upper bounded by an implicit alignment loss and the reconstruction target aligns the features of the ³⁸ masked view and the unmasked view of the same image. Intuitively, some masked views of different

images in the same class can be very similar, which means different images can be bridged by these
 similar overlapped views.

With the analysis of the MAE loss, we propose a guarantee for MAE training process and it will
converge with some practical assumptions. Based on our analysis of MAE training process, we
propose a Uniformity-promoting MAE (U-MAE) loss, which adds a regularizer term to the original
MAE loss. Empirically, we find our proposed loss can significantly improve the performance of MAE
on different benchmark datasets (CIFAR-10, ImageNet-100) with different backbones (ViT-Base,
ViT-Large). We summarize our contributions as follows:

- We prove that MAE loss is upper bounded by an implicit alignment loss and develop a new understanding of how MAE bridges the samples in the same class. Based on our analysis, we propose a new U-MAE loss.
- We establish a guarantee for the downstream performance of MAE with practical assumptions and compare the training process of contrastive learning and MAE.

• We empirically verify that our proposed U-MAE loss can significantly improve the performance of MAE across different real-world datasets, including CIFAR-10 and ImaegNet-100.

54 2 Related work

55 **Self-Supervised Learning.** With meaningful surrogate objectives, self-supervised learning has 56 gradually closed the performance gap between supervised and unsupervised learning without access to expensive supervised information. For example, contrastive learning designs appropriate data 57 augmentations to draw together the semantically similar samples and obtain impressive downstream 58 performance on various datasets [3, 10, 8]. Masked Image Modeling (MIM) is another surrogate 59 task to learn meaning representations, which masks some patches of samples and trains the encoder 60 based on the comparison between the masked and unmasked patches. Masked Autoencoder (MAE) 61 [11] is a representative work of MIM and obtains state-of-the-art performance on various large-scale 62 datasets like ImageNet. The core task of MAE is to reconstruct masked images with an asymmetric 63 Vision Transformer architecture 64

The Theory of Self-supervised Learning. Inspired by the empirical success of self-supervised 65 learning, many researches try to establish theoretical analysis on self-supervised learning. [15] 66 establishes the bound between supervised learning and unsupervised learning. [16] understands 67 contrastive learning from a mutual information perspective. Analyzing self-supervised learning from 68 a graph perspective has been proved to be feasible in contrastive learning. [9] analyzes contrastive 69 learning with the spectral graph theory and gives a bound for the downstream performance related 70 to the spectral properties of the augmentation graph. [18] constructs an augmentation graph and 71 analyzes the gap between contrastive learning and supervised learning with it. In this paper, we also 72 try to understand the training process of MAE from the graph perspective, as we find that contrastive 73

r4 learning and MAE share a significant similarity, *i.e.*, the samples in the same class could generate quite similar views with transformations, including data augmentations used in contrastive learning and mask used in MAE. However, the analysis in this paper is not a trivial solution and there exist some significant differences between contrastive learning and MAE. For example, the existing theory can't explain the relationship between the mask ratio and downstream performance, and the graph of

⁷⁹ MAE is a directed graph instead of an undirected graph as analyzed in contrastive learning.

Understanding MAE Training Process. Despite the impressive success of MAE, the theoretical 80 understanding of it is still underexplored. [2] is the first work to understand MAE through a 81 mathematical viewpoint, which explains the success of MAE from a perspective of an integral kernel. 82 They analyze the optimizing process of MAE and show the importance of patchifying and the decoder. 83 However, some designs and properties of MAE are still not fully understood. For example, why MAE 84 can learn meaningful representation for downstream tasks with an instance-level reconstruction task 85 and why the mask ratio needs to be set to 0.75? In this paper, we try to further explain this powerful 86 self-supervised training framework from a different perspective, *i.e.*, graph perspective, and answer 87 these questions. 88

3 Masked AutoEncoding as Implicit Contrastive Learning

90 3.1 Problem Formulation

We begin by introducing the basic notations and common process of MAE, which consists of two stages, *i.e.*, MAE pretraining and supervised finetuning. In the first stage, we pretrain an encoderdecoder architecture to reconstruct the masked images. In the second stage, we evaluate the quality of the learned representations with labeled dataset $\mathcal{D}_l = \{(x_i, y_i) | x_i \in \mathbb{R}^d, 0 \le y \le k\}$ with a classification task.

96 **3.1.1 MAE Pretraining**

Assuming access to an unlabeled dataset of natural images $\mathcal{D}_u = \{\bar{x}_i | \bar{x}_i \in \mathbb{R}^d\}$, we denote the natural data distribution as $P_d(\bar{x})$. Masked Autoencoders (MAE) [11] learn data representations with the following procedures.

Pretext Generation - Random Patch-wise Masking. First, for a random image $\bar{x} \in \mathcal{D}_u$, it could be (equivalently) patchified as a matrix, denoted as $\hat{x} \in \mathbb{R}^{n \times p}$, where *n* denotes the number of patches, and *p* denotes the patch size, *e.g.*, 4×4 . Second, according to a certain masking ratio $\rho \in [0, 1]$, we randomly draw a binary mask $m \in \{0, 1\}^n$ for each patch. Denote the number of non-zero elements of *m* as |m|. and selecting patches according to the mask *m* leads to two complementary views,

$$x_1 = \hat{\bar{x}}[m] \in \mathbb{R}^{|m| \times p}, \quad x_2 = \hat{\bar{x}}[1-m] \in \mathbb{R}^{(n-|m|) \times p},$$
 (1)

where $[\cdot]$ denote the index selection function. Overall, we denote this random masking process for generating x_1, x_2 from x as $M(x_1, x_2 | \bar{x})$, and their marginal distribution as $M_1(x_1 | \bar{x})$ and $M_2(x_2 | \bar{x})$, respectively. Here $\mathcal{M}, \mathcal{M}_1, \mathcal{M}_2$ refer to the probability density functions of the corresponding distributions. As long as $\rho \neq 0.5$, the two distributions are different, *i.e.*, $\exists x, s.t. M_1(x | \bar{x}) \neq$ $M_2(x | \bar{x})$. And we denote $p_x \in \mathbb{R}^{n_x \times s}$ as the position information of the view $x \in \mathbb{R}^{n_x \times p}$ with n_x patches, typically represented as Position Encodings (PE).

Reconstruction-based Pretraining. According to He *et al.* [11], the MAE framework is as follows: an encoder f first maps (x_1, p_{x_1}) to a latent feature $z = f(x_1, p_{x_1}) \in \mathbb{R}^{|m| \times q}$. Afterwards, a decoder g takes z_1 and p_{x_2} as the input, and maps to the sample space to get an estimation of x_2 , *i.e.*, $\hat{x}_2 = g(z_1, p_{x_2}) = g(f(x_1, p_{x_1}), p_{x_2}) \in \mathbb{R}^{(n-|m|) \times p}$. At last, the encoder-decoder architecture of MAE is learned with a simple ℓ_2 reconstruction loss,

$$\min_{f,g} L_{rec}(\mathcal{D}_u; f, g); \quad L_{rec}(\mathcal{D}_u; f, g) = \mathbb{E}_{P_d(x)} \mathbb{E}_{M(x_1, x_2|x)} \left\| g(f(x_1, p_{x_1}), p_{x_2}) - x_2 \right\|^2.$$
(2)

116 **3.1.2 Linear Evaluation**

After pretraining, we can evaluate the quality of the extracted data representations via a linear classification task. Specifically, assuming access to the labeled dataset $\mathcal{D}_l = \{(x_i, y_i) | x_i \in \mathbb{R}^d, y \in \mathbb{R}^d, y \in \mathbb{R}^d\}$ 119 [k], we train a linear classifier $l : \mathbb{R}^q \to \mathbb{R}^k$ with the encoder held fixed

$$\min_{l} L_{cls}(\mathcal{D}_l; f, l); \quad L_{cls}(\mathcal{D}_l; f, l) = \mathbb{E}_{P_d(x, y)} \ell_{ce}(l(f(x)), y), \tag{3}$$

where ℓ_{ce} denotes the cross entropy loss.

121 3.2 MAE Loss is upper bounded by an Implicit Alignment Loss

Reconstruction Loss is upper bounded by an Implicit Alignment Loss. We begin with an analysis on the reconstruction loss used in the training process of MAE. From its form, we can find that it is a common autoencoding loss that tries to reconstruct samples with an encoder-decoder architecture. As MAE uses the transformer network as the encoder and the decoder, so we can obtain some properties to simplify our subsequent analysis on the training process of MAE. [20] analyzes the universal approximation of the transformer and proposes the following results,

Lemma 3.1 (Theorem 3 [20]). We denote F_{CD} as the set of all continuous functions that map a compact domain. Let $1 \le p < \infty$ and $\varepsilon > 0$, then for any given $f \in F_{CD}$, there exists a Transformer network $g \in T_P^{2,1,4}$ such that we have $d_p(f,g) \le \varepsilon$.

Combined with the universal approximation property of the transformer network and the encoderdecoder architecture, we propose the following assumption which states there exists a pseudo-inverse encoder of the decoder, which maps from the complementary view (x_2, p_{x_2}) back to the representation

134 space of the encoder.

Assumption 3.2. For the decoder function $g : \mathbb{R}^{|m| \times q} \times \mathbb{R}^{(n-|m|) \times s} \to \mathbb{R}^{(n-|m|) \times p}$, we denote the pseudo-inverse encoder of it as $f_g : \mathbb{R}^{(n-|m|) \times p} \times \mathbb{R}^{(n-|m|) \times s} \to \mathbb{R}^{|m| \times q}$, satisfying

$$\forall (x, p_x) \in \mathcal{X} \times \mathcal{P}_{\mathcal{X}}, \quad \|g(f_g(x, p_x), p_x) - x\|^2 \le \varepsilon.$$
(4)

137 138 Assumption 3.2 characterizes the reconstruction ability of an encoder-decoder architecture. As it does not involve the mask technique, it amounts to a simple autoencoding model, whose reconstruction 139 error can be bounded theoretically and empirically with powerful neural networks [20, 13]. Thus, 140 this assumption is likely to hold. However, we notice that a trivial autoencoding loss can not learn 141 meaningful representations for downstream performance. To verify that MAE reconstruction loss is 142 not only a simple autoencoding loss, we set the mask ratio to 0 and find that the linear accuracy of 143 learned representation on ImageNet-100 decreases to 13.8%. Then we further explore the properties 144 of MAE loss with the above practical assumption, 145

$$L_{rec}(f,g) = \mathbb{E}_{x_1,x_2} \|g(f(x_1, p_{x_1}), p_{x_2}) - x_2\|^2$$

=\mathbb{E}_{x_1,x_2} \|g(f(x_1, p_{x_1}), p_{x_2}) - g(f_g(x_2, p_{x_2}), p_{x_2}) + g(f_g(x_2, p_{x_2}), p_{x_2}) - x_2\|^2 \quad (5)
$$\leq 2\mathbb{E}_{x_1,x_2} \|g(f(x_1, p_{x_1}), p_{x_2}) - g(f_g(x_2, p_{x_2}), p_{x_2})\|^2 + 2\varepsilon.$$

From Equation (5), we find the reconstruction loss used in MAE is upper bounded by an implicit 146 alignment loss which pulls in the distance of different masked views encoded by the encoder-decoder 147 architecture. Note that MAE only uses the encoder for downstream tasks, so then we consider the 148 smoothness of the decoder to analyze whether reconstruction loss aligns the different views in the 149 encoder feature space. [12] analyzes the Lipschitz continuity of a transformer network and prove 150 that a transformer network is Lipschitz continuous when we give constraints to the architecture. We 151 also empirically analyze the Lipschitz constant in the input domain of the decoder during the training 152 $\frac{\|g(f(x_1,p_{x_1}),p_{x_2}) - g(f_g(x_2,p_{x_2}),p_{x_2})\|^2}{\|f(x_1,p_{x_1}) - f_g(x_2,p_{x_2})\|^2}$ and we find it keepes smaller process of MAE, *i.e.*, $\min_{x_1, x_2 \in \mathcal{D}_u}$ 153 than 1.42. More details about the empirical verification can be found in Appendix B.1. Combined 154 with the empirical verification, it's practical to assume that the decoder q is L-Lipschitz, *i.e.*, 155

156 Assumption 3.3. For the decoder function g, we assume there exists a constant L, satisfying

$$\forall (x_1, x_2), \quad \|g(x_1) - g(x_2)\|^2 \le L \|x_1 - x_2\|^2.$$
(6)

Theorem 3.4. Under Assumptions 3.2 & 3.3, the MAE reconstruction loss can be upper bounded by the alignment loss between \tilde{x}_1 and \tilde{x}_2 , where $\tilde{x} = [x, p_x]$ denotes the position-augmented input:

$$L_{rec}(f,g) \le -2L \cdot \mathbb{E}_{\tilde{x}_1,\tilde{x}_2} \|f(\tilde{x}_1)^\top f_g(\tilde{x}_2)\|^2 + 2\varepsilon + C.$$

$$\tag{7}$$

160 Here ε is the approximation error of f_q and C is a constant.



Figure 2: Appropriate mask ratio can generate similar views from different samples in the same class.

161 Intuitively, Theorem 3.4 shows that the reconstruction loss of MAE implicitly minimizes the distance

of different masked views of the same sample in the encoder features space. It can be seen that this is similar to the alignment objective that appears in contrastive (as well as non-contrastive) methods, *i.e.*, $L(f,g) = -2\mathbb{E}_{x_1,x_2}f(x_1)^{\top}g(x_2)$. However, there exist obvious differences as the views used in MAE is asymmetric with different mask ratio while the views in canonical contrastive learning are symmetric with data augmentations. Meanwhile, MAE could avoid latent collapse because f_g is always a nontrivial target.

Suggested Uniformity-promoting MAE Loss. With the proof above, we can find that the MAE loss is upper bounded by an implicit alignment loss, which aligns features of the masked and unmasked views. [17] understands the alignment loss from two properties, *i.e.*, alignment and uniformity. As proposed in [17], without keeping the uniformity of the features, the cluster performance of the same class will be significantly hurt. So we add a regularizer term to ensure the uniformity of the encoder to improve the performance of MAE and propose our suggested Uniformity-promoting MAE (U-MAE) loss,

$$L_{\text{U-MAE}}(f,g) = \mathbb{E}_{\tilde{x}_1,\tilde{x}_2} \|g(f(x_1, p_{x_1}), p_{x_2}) - x_2\|^2 - \lambda \cdot \mathbb{E}_{\tilde{x}_1,\tilde{x}_2^-} \|f(\tilde{x}_1) - f_g(\tilde{x}_2^-)\|^2, \quad (8)$$

where \tilde{x}_1, \tilde{x}_2 is the unmasked and masked views of the same sample, *i.e.*, the positive samples, and \tilde{x}_2^- is the masked view of an independently drawn sample, *i.e.*, the negative sample. As the implicit reverse encoder f_g is a hypothetical function and hard to obtain, we replace it with the encoder f, and obtain the following practical U-MAE loss with symmetric uniformity regularizer:

$$L_{\text{U-MAE}}^{\text{sym}}(f,g) = \mathbb{E}_{\tilde{x}_1,\tilde{x}_2} \|g(f(x_1, p_{x_1}), p_{x_2}) - x_2\|^2 - \lambda \cdot \mathbb{E}_{\tilde{x}_1,\tilde{x}_1^-} \|f(\tilde{x}_1) - f(\tilde{x}_1^-)\|^2, \quad (9)$$

where \tilde{x}_1^- is an independently drawn unmasked view. Intuitively, a large distance between encoder outputs can also effectively promote feature diversity.

181 4 Generalization Theory

182 4.1 Spectral Formulation

Based on our above analysis of MAE loss, we know that MAE loss is upper bounded by an implicit alignment loss. In this section, we revisit the training process of MAE from a spectral graph perspective and find that MAE reconstruction loss is closely related to an asymmetric matrix factorization problem. Built upon this connection, we establish theoretical guarantees on the downstream performance and characterize its influencing factors.

Constructing the Directed Augmentation Graph. Following [9], we define the population augmentation graph $G(\mathcal{X}, A)$, where \mathcal{X} denotes all masked views of size $N = |\mathcal{X}|$, and $A = (w_{xx'}) \in \mathbb{R}^{N \times N}$ denotes the adjacency matrix filling with edge weight $w_{xx'}$ for two masked views $x, x' \in \mathcal{X}$. Specifically, we define the weight $w_{xx'}$ as the marginal probability of generating the pair x and x'from the random natural data

$$w_{xx'} = \mathbb{E}_{\bar{x}\sim\mathcal{P}_d}[M_1(x|\bar{x})M_2(x'|\bar{x})].$$
(10)

Intuitively, the population augmentation graph formalizes the process of how mask technique bridges different samples in the same ground-truth class. For example, the tires of two different cars could be quite similar and can be seen as generated from the same sample, then some masked views of these two cars will have large weights between them. Obviously, the samples of the same class will have larger weights then the samples of different classes. Note that different from the undirected augmentation graph considered in [9], the graph $G(\mathcal{X}, w)$ is directed, *i.e.*,

$$\mathbb{E}_{\bar{x}\sim\mathcal{P}_d}[A_1(x|\bar{x})A_2(x'|\bar{x})] = w_{xx'} \neq w_{x'x} = \mathbb{E}_{\bar{x}\sim\mathcal{P}_d}[A_1(x'|\bar{x})A_2(x|\bar{x})].$$
(11)

Accordingly, we can normalize the directed adjacency matrix A as $\bar{A} := D_{row}^{-1/2} A D_{col}^{-1/2}$, where D_{row}, D_{col} are diagonal matrices with $\forall x \in \mathcal{X}, (D_{row})_{xx} = w_{x} := \sum_{x'} w_{xx'}, (D_{col})_{xx} =$ $w_{\cdot x} := \sum_{x'} w_{x'x}$. Because \bar{A} is asymmetric, generally $w_{x} \neq w_{\cdot x}$.

Asymmetric Matrix Factorization. With the directed augmentation graph, we revisit our U-MAE
 loss and find it can be upper bounded by an Asymmetric Matrix Factorization loss.

Theorem 4.1. For the augmentation graph adjacency matrix A, we consider two matrices $U \in \mathbb{R}^{n \times k}$ and $V \in \mathbb{R}^{k \times n}$, and a decomposition objective $L_{mf}(U, V) = \|\bar{A} - UV^{\top}\|_{F}^{2}$. We have

$$L_{U-MAE}(f,g) \le 2L_{mf}(f,f_g) + 2\varepsilon + C, \tag{12}$$

where ε is the approximation error of the transformer network.

Built upon this connection, in the next step, we will discuss the downstream performance by analyzing the asymmetric matrix decomposition objective. To achieve this, we first introduce two common assumptions adopted in previous work [9, 18], one on the label consistency of masking and one on the expressivity of the hypothesis class. Both are likely to hold in practice as long as we choose proper mask ratio and architecture for MAE.

Assumption 4.2 (Labels are recoverable from masked views). Let $\bar{x} \in D_u$ be a sample and $y(\bar{x})$ be its label. For a masked view x_1 of x, we assume there exists a classifier f that can predict $y(\bar{x})$ given x with error at most α . That is, $f(x_1) = y(\bar{x})$ with probability at least $1 - \alpha$.

Assumption 4.3 (Expressivity of the hypothesis class). Let \mathcal{F}, \mathcal{G} be the hypothesis classes of the encoder f and the decoder g, respectively. We assume that at least one of the global minima of $L_{U-MAE}(f, g)$ belongs to \mathcal{F} and \mathcal{G} .

In the following, combined with our analysis on Theorem 4.1, we establish a bound between the downstream performance of MAE and the spectral property of the augmentation graph. The details of the proof can be found in Appendix A.3.

Theorem 4.4. Under Assumptions 3.2, 3.3, 4.2 & 4.3, let \mathcal{L} be the Laplacian matrix of the adjacent matrix A and $\lambda_1, \lambda_2, \dots, \lambda_k, \lambda_{k+1}$ be be the k+1 smallest eigenvalues of \mathcal{L} . $\zeta(f_{U-MAE}^*)$ is denoted as the downstream error for the minimizer of MAE loss f_{U-MAE}^* , we have

$$\zeta(f_{U-MAE}^{\star}) \le O(L\sqrt{2\alpha}/\lambda_{k+1}) + 2\varepsilon.$$
(13)

224 225 From Theorem 4.4, we can conclude that the downstream performance can be minimized with a small Lipschitz constant L, a small decoder inversion error ε , a small the labeling error α , and a large 226 k + 1-th eigenvalue of the augmentation graph λ_{k+1} . In other words, MAE with better landscape 227 smoothness and invertible decoder will be likely to have a smaller downstream error, and in the 228 meantime the masking ratio ρ should be properly chosen such that the labeling error α is small and 229 λ_{k+1} is large. The last condition seems less intuitive. In fact, according to spectral graph theory [4], 230 a large eigenvalue λ_{k+1} stems from a strong connectivity of the augmentation graph \mathcal{G} . Thus, MAE 231 often requires a large masking ratio to bridge different samples together. Meanwhile, as illustrated 232 in Figure 2, a too large mask ratio, e.g., $\rho = 0.95$, could make it hard to recover the label from 233 the masked views, which results in a large α . As a result, Theorem 4.4 suggests that there exists a 234 tradeoff in the mask ratio ρ , and we should choose a properly large one, e.g., $\rho = 0.75$. 235

236 4.2 Empirical Investigation of the Directed Augmentation Graph of MAE

Properties. As discussed in [18], contrastive learning bridges the samples in the same class by the similar views generated by appropriate data augmentations, comparing to the augmentation graph



(a) The distance between (b) The distance between (c) The distance between (d) The relative distance intra-class samples inter-class samples overall samples

between intra-class samples and inter-class samples

Figure 3: The influence of mask ratio on the average L2 distance on ImageNet-100 between (a) intra-class samples (b) The inter-class samples (c) intra-class samples, and the relative distance between intra and inter-class samples.

used in MAE, we can find that both of them align the intra-class samples by transforming different 239 anchor samples into similar views. However, although they have a lot in common, there exist obvious 240 differences between contrastive learning and MAE. Firstly, as the mask ratio is not equal to 0.5, 241 the two views of an anchor sample in MAE are asymmetric, which means the augmentation graph 242 changes from an undirected graph in contrastive learning into a directed graph in MAE. So we can not 243 directly use the theoretical framework designed for analyzing the undirected augmentation graph in 244 contrastive learning and we find it can be analyzed with an Asymmetric Matrix Factorization objective. 245 And as the mask ratio of similar views of different anchors should be equal, the augmentation graph 246 of MAE becomes a bipartite graph, which means views of different samples with the same mask 247 ratio can't be connected directly and they need a middleman between them. What's more, the two 248 views of the same anchor in contrastive learning is independent while the two views in MAE loss are 249 dependent and contain the whole anchor, which means the two views in MAE training can represent 250 most of the features of the sample. 251

The Necessity of High Mask Ratio. There exists a difficulty in understanding why the training 252 process of MAE needs a high mask ratio. We find that a major reason is that the higher mask ratio 253 will generate more similar views of intra-class samples. For example, in Figure 2, we can find that 254 when the mask ratio increases to 75%, the views of two different cars can be concentrated on the 255 features of similar tires. To verify our perspective, we conduct an experiment on images of CIFAR-10 256 with different masked ratios. We evaluate the distance between two images by calculating the l_2 257 distance between the patches of them. As shown in Figure 3(a), we can find that the average distance 258 of intra-class images decreases with the enhancement of mask ratio, which means a higher mask ratio 259 260 leads to generating more edges between intra-class samples.

Drawbacks of Too High Mask Ratio. On the other hand, as the highest mask ratio can not lead 261 to the highest downstream performance of MAE, we find it is related to the decreased size of the 262 support set of an anchor, *i.e.*, some transformed views of an anchor will be very similar. Indeed, 263 as shown in Figures 3(b) & 3(c), the average L_2 distance between inter-class samples and overall 264 265 samples decreases as the mask ratio increases, which means when we take an excessive mask ratio, numerous views will be merged and the support set of the anchors will become smaller. And the 266 views with an extremely high mask ratio will lose the feature of their class. For example, as shown in 267 Figure 2, we can find that when the mask ratio arrives at 95%, the features of the car disappear. 268

The Sweet Spot for Mask Ratio. The discussion above reveals that either too small or too high 269 mask ratios can result in bad features, and a good mask ratio should be selected such that intra-class 270 distance is relatively high while inter-class distance is relatively low. Motivated by this, we plot the 271 relative distance (intra-class over inter-class) in Figure 3(d), and we can find that the relative distance 272 decreases first with larger mask ratio, showing that the intra-class distance decreases faster than the 273 inter-class distance under small mask ratio. When $\rho > 0.7$, the relative distance becomes larger 274 again, indicating that the difference between intra-class and inter-class edges disappears under too 275 large mask ratio. The sweet spot lies in $\rho = 0.7$, which is pretty close to the optimal masking rate of 276 MAE is $\rho = 0.75$ [11]. This shows that our theoretical analysis of the effect of MAE ratio agrees 277 surprisingly well with the practice of MAE. 278

Table 1: Linear evaluation accuracy and finetune accuracy on CIFAR-10 and ImageNet-100 pretrained with MAE loss and U-MAE loss with different ViT backbones. The uniformity regularizer term significantly improves the linear evaluation performance of MAE loss with different backbones across different real-world datasets without hurting the performance of finetune accuracy.

	Objective	CIFAR-10		ImageNet-100	
BackBone		ViT-Tiny	ViT-Base	ViT-Base	ViT-Large
Linear Evaluation	MAE	52.0	59.9	37.5	39.5
	U-MAE	69.4	72.0	56.3	61.4
Finetune Results	MAE	89.6	90.7	86.9	87.3
	U-MAE	89.4	90.8	86.8	87.3

279 5 Experiments

280 5.1 Evaluation on Benchmark Datasets

We conduct experiments of our U-MAE loss on various real-world datasets, including CIFAR-10 [14] and ImageNet-100 [5]. In the following, we introduce our empirical setup on CIFAR-10 and ImageNet-100.

Setup. As our proposed loss is a promoting MAE loss, we mainly follow the basic setup of MAE. 284 We use ViT as the encoder and use the flexible decoder as proposed in MAE. To further present the 285 performance of our loss, we carry out our experiment on different variants of ViT, *i.e.*, ViT-Tiny, 286 ViT-Base and ViT-Large. We follow the default setting as proposed in [11]. We use the recommended 287 mask ratio, 75%. For the uniformity term of our proposed loss, we set the coefficient of the uniformity 288 term to 0.001. For ImageNet-100, we pretrain the model for 200 epochs with batch size 256 and 289 weight decay 0.05. And for CIFAR-10, we pretrain the model for 2000 epochs with batch size 4096 290 and weight decay 0.05. On the stage of supervised finetuning, we train a supervised classifier on the 291 pretrained encoder. We conduct both linear evaluation and non-linear finetuning on the unsupervised 292 pretrained encoder. For linear evaluation, we train a linear classifier on the frozen pretrained encoder. 293 As for non-linear finetuning, we train both the pretrained encoder and the linear classifier with Cross 294 Entropy loss. 295

Results. In Table 1, we compare the linear evaluation results between original MAE loss and U-296 MAE loss on different benchmarks, including CIFAR-10 and ImageNet-100, and different backbones, 297 including ViT-Tiny, ViT-Base and ViT-Large. We find that our promoting loss increases 14.71% 298 for linear evaluation results on CIFAR-10 with two different backbones and increases 20.35 % on 299 ImageNet-100 with two different backbones. And in Table 1, we also present the comparison of 300 finetuning results. We find that our proposed loss will not hurt the performance of finetuning results 301 of MAE. Our results empirically verify the effectiveness of our loss and show that the U-MAE 302 303 loss achieves better performance than the original MAE loss across different datasets and different backbones. 304

305 5.2 Empirical Understanding and Ablation Study.

Visualization of Representations. To intuitively understand the improvement of our U-MAE 306 loss on clustering intra-class samples, we use t-SNE to visualize the representations trained with 307 MAE loss and our proposed loss on ten random class of ImageNet-100 datasets detailed classes are 308 introduced in Appendix B.4. We find that with our uniformity regularizer term, the samples are much 309 better-clustered corresponding to their ground-truth labels. To be specific, the red class (representing 310 for hens) and the gray class (representing for indigo birds) are separated from others. We note most 311 of other classes are the animals living in the oceans) while these two types are more like the birds 312 (living on the land or the sky). So these two classes are more easier to be distinguished, especially 313 with our uniformity regularizer term. 314

Different Coefficients of the Regularizer Term. The most important hyper-parameter of our proposed loss is the coefficient of the uniformity regularizer term. In Figure 5, we present the results of linear evaluation on ImageNet-100 trained with the U-MAE loss with different coefficients of the



Figure 4: Visualization of representations on random 10 classes of ImageNet (0-9 classes described in Appendix B.4) trained with MAE loss and our U-MAE loss. Our loss significantly improves the class-clustering performance of the encoder.



(a) Different coefficient of the uniformity (b) Comparison between original MAE regularizer term. and our U-MAE along training.

Figure 5: (a) The linear evaluation results trained with different coefficients of the uniformity regularizer term of our U-MAE loss. It shows that an appropriate coefficient is important for our proposed loss. (b) The linear evaluation results during the training process. It shows that our proposed loss improve the downstream performance of MAE with different training time.

regularizer term. We can find that when the coefficient is 0, our proposed loss is an original MAE loss. The downstream performance increases when the coefficient increases from 0 to 0.001. However, the overlarge coefficient will also hurt the performance of our proposed loss as the task of MAE will be overlooked.

Training process. To further compare the performance between the original MAE loss and our U-MAE loss, we report the linear evaluation accuracy on ImageNet-100 during the training process. In Figure 5, we show the linear evaluation results during the training process. We can observe that our proposed loss improves the performance of MAE with all different training epochs, which verifies that our proposed loss is a promoting loss of MAE for both short and long time learning.

327 6 Conclusion

In this paper, we propose a new theoretical understanding of MAE. With the analysis of MAE loss, we find that MAE loss is upper bounded by an implicit alignment loss. Then we explain the training process of MAE from a graph perspective and establish a guarantee for its downstream performance. Based on our theory, we propose an Uniformity-promoting MAE (U-MAE) loss and verify that it can significantly improve the downstream performance of MAE across different benchmark datasets, including CIFAR-10 and ImageNet-100.

334 **References**

- [1] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. BEit: BERT pre-training of image transformers. In *ICLR*, 2022.
- [2] Shuhao Cao, Peng Xu, and David A. Clifton. How to understand masked autoencoders. *CoRR*, 2022.
- [3] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework
 for contrastive learning of visual representations. In *ICML*, 2020.
- [4] Fan Chung. Laplacians and the cheeger inequality for directed graphs. *Annals of Combinatorics*, 9(1):1–19, 2005.
- [5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
 hierarchical image database. In *CVPR*, 2009.
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018.
- [7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,
 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image
 recognition at scale, 2020.
- [8] Jean-Bastien Grill, Florian Strub, Florent Altché, C. Tallec, Pierre H. Richemond, Elena
 Buchatskaya, C. Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi
 Azar, B. Piot, K. Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent: a
 new approach to self-supervised learning. In *NeurIPS*, 2020.
- [9] Jeff Z HaoChen, Colin Wei, Adrien Gaidon, and Tengyu Ma. Provable guarantees for self supervised deep learning with spectral contrastive loss. In *NeurIPS*, 2021.
- [10] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for
 unsupervised visual representation learning. In *CVPR*, 2020.
- [11] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked
 autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377*, 2021.
- [12] Hyunjik Kim, George Papamakarios, and Andriy Mnih. The lipschitz constant of self-attention.
 In *ICML*, 2021.
- ³⁶³ [13] Diederik P Kingma and Max Welling. Auto-encoding variational Bayes. In *ICLR*, 2014.
- [14] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
 2009.
- [15] Nikunj Saunshi, Orestis Plevrakis, Sanjeev Arora, Mikhail Khodak, and Hrishikesh Khande parkar. A theoretical analysis of contrastive unsupervised representation learning. In *ICML*,
 2019.
- [16] Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola. What
 makes for good views for contrastive learning. In *NeurIPS*, 2020.
- [17] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through
 alignment and uniformity on the hypersphere. In *ICML*, 2020.
- [18] Yifei Wang, Qi Zhang, Yisen Wang, Jiansheng Yang, and Zhouchen Lin. Chaos is a ladder: A
 new theoretical understanding of contrastive learning via augmentation overlap. In *ICLR*, 2022.
- [19] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han
 Hu. Simmim: A simple framework for masked image modeling. In *CVPR*, pages 9653–9663,
 2022.
- [20] Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar.
 Are transformers universal approximators of sequence-to-sequence functions? In *ICLR*, 2020.

[21] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong.
 Image BERT pre-training with online tokenizer. In *ICLR*, 2022.

382 Checklist

383	1. For all authors
384 385	 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
386	(b) Did you describe the limitations of your work? [No]
387	(c) Did you discuss any potential negative societal impacts of your work? [No]
388 389	(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
390	2. If you are including theoretical results
391	(a) Did you state the full set of assumptions of all theoretical results? [Yes]
392	(b) Did you include complete proofs of all theoretical results? [Yes]
393	3. If you ran experiments
394 395	(a) Did you include the code, data, and instructions needed to reproduce the main experi- mental results (either in the supplemental material or as a URL)? [No] Release upon
396	acceptance.
397 398	(b) Did you specify all the training details (e.g., data splits, hyperparameters, now they were chosen)? [Yes] See Section 5.1
399 400	(c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [No]
401 402	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [No]
403	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
404	(a) If your work uses existing assets, did you cite the creators? [Yes]
405	(b) Did you mention the license of the assets? [No]
406	(c) Did you include any new assets either in the supplemental material or as a URL? [No]
407	(d) Did you discuss whether and how consent was obtained from people whose data you're
408	using/curating? [No]
409	(e) Did you discuss whether the data you are using/curating contains personally identifiable
410	information or offensive content? [No]
411	5. If you used crowdsourcing or conducted research with human subjects
412 413	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
414 415	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
416 417	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]