
Mask Models are Token Level Contrastive Learners

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In recent years, the field of self-supervised learning has seen a surge in the develop-
2 ment of mask models, which have been demonstrated to have strong performance
3 on downstream tasks and efficient training. To better understand the underlying
4 mechanism behind these models' success, we propose a theoretical framework for
5 understanding mask models. By treating mask modeling as a low-rank recovery
6 task, we demonstrate that it is a parametric version of Spectral Clustering and
7 the reconstruction loss conforms to the form of Spectral Contrastive loss. This
8 means that mask modeling can be understood as a token level Contrastive Learning.
9 Our framework can be used to explain why optimal masking ratios vary among
10 modalities and why there is a large gap between linear probing and finetuning
11 performance for mask models. Additionally, our analysis suggests that the success
12 of mask models depends on the model architecture, where a token mixing layer
13 and layer normalization are crucial for the success of mask models. Our framework
14 has the potential to be a step stone for future algorithm and network architecture
15 design in the field of self-supervised learning.

16 1 Introduction

17 With the rapid-growth of deep learning and its increasing demand for data, self-supervised learning
18 arises as a research topic in-demand. Among the successful self-supervised learning models, mask
19 models have received significant attention for their strong downstream performance and efficient
20 training [15, 38, 12, 5, 23, 34, 18, 4]. However, mask-modeling has long been regarded as an
21 engineering trick, and its underlying mechanism remains poorly understood.

22 Empirically, we have observed that Mask Image Modelling (MIM) leads to varying performance
23 improvements on different downstream tasks compared to previous baselines [15]. Specifically,
24 MIM has been found to perform better on fine-grained tasks such as semantic segmentation and
25 object detection, compared to classification tasks. This phenomenon leads us to hypothesize that
26 the representation learned by MIM models is fundamentally similar to that of image segmentation
27 (clustering).

28 In this work, we propose a theoretical analysis of mask modelling by treating it as a low-rank recovery
29 (LRR) task. Our analysis further demonstrates that the reconstruction loss can be rewritten as a
30 Contrastive loss [10].

31 The LRR problem aims to find the low-rank approximation of a given matrix, and has been used as a
32 method for subspace clustering [28]. Additionally, as the optimal solution of the LRR problem is a
33 combination of leading eigenvectors, we are naturally led to Spectral Clustering, which also utilizes
34 leading eigenvectors [32, 26]. Our results show that MIM approximates the Spectral Clustering
35 features of an image-related graph, where each node represents a patch of the image.

36 By viewing the Masked Image Model (MIM) as a parametric version of Spectral Clustering, we can
37 rewrite the reconstruction loss of mask models in the form of Spectral Contrastive loss on the token

38 level [14]. This allows MIM to be viewed as a token-wise Contrastive Learning method, which
39 attracts similar patches while repelling dissimilar ones, resulting in smaller distances within clusters
40 and larger distances between clusters. However, there are some key differences between mask models
41 and traditional Contrastive Learning methods. Specifically, mask models operate on the token level,
42 whereas traditional Contrastive Learning methods focus on the global feature of the entire input, and
43 in mask models, positive samples are not clearly defined, but are "randomly sampled" based on the
44 similarity between Spectral Clustering features.

45 Based on the formulation, we could answer several concerning questions about mask models: 1) Why
46 optimal masking ratio vary among modalities? 2) Why is there a large gap between linear probing and
47 finetuning performance for mask models? 3) Does mask modelling rely on network architectures?

48 For the first question, we argue that a critical factor that affects the goodness of pretrained features is
49 the number of clusters in Spectral Clustering. For example, if we have an image with a dog on the
50 grass, intuitively we should have two clusters: grass and dog. It could be less representative if we
51 have more clusters and divide one of the existing clusters into different sub-clusters and repel one
52 from each other. The number of clusters is given by the number of leading eigenvectors, which is
53 related to the rank of reconstructed matrix in LRR problem and masking ratio in MIMs. This explains
54 why we need different masking ratios in different modalities [12, 15, 34, 18].

55 For the second question, it is due to the nature of token level Contrastive learning. Pretrained mask
56 models learn to divide tokens into clusters, but doesn't always learn which cluster is most related
57 to the class. Therefore, token mixing layers are needed to "select" clusters. Most MIMs apply an
58 extra BatchNorm layer when performing linear probing, otherwise a huge accuracy drop is witnessed
59 [19, 15]. It could be due to the lack of patch selection and a BatchNorm is needed to add non-linearity.
60 In contrast, we found that partially finetuning one linear layer for row mixing with the prediction
61 head could much improve classification accuracy.

62 For the third question, the answer is "Yes". Model architectures containing token mixing layers
63 plays a crucial role in the success of mask models in classification tasks [35, 13, 24, 33]. Finetuning
64 these layers allows the model to learn how to select tokens. Meanwhile, the layer normalization
65 in the decoder might also be important, as it serves as a token level batch normalization, which is
66 commonly used in the projection layer of Contrastive Learning models to improve performance
67 [3, 19]. Therefore, we conclude that mask model is dependant of network architecture.

68 In a summary, our main contributions are:

- 69 1. We created a mathematical framework for mask image modeling by viewing it as a low-rank
70 recovery problem.
- 71 2. We found that mask model could be viewed as a token level Contrastive Learning, which
72 could account for its good performance on downstream tasks.
- 73 3. Our analysis framework could explain several important behaviors of mask models and
74 guide future model architecture design and parameter choosing.

75 We mainly conducted experiments on images, but our findings could be easily generalized to all
76 modalities.

77 **2 Related Works**

78 **2.1 Mask Image Modelling**

79 The recent trend in self-supervised learning is to train vision transformers using masked images to
80 reconstruct the original ones. [13]. Different types of reconstruction objectives, such as token-wise,
81 feature-wise, and pixel-wise reconstruction, are being tested. [39, 8, 7]. These kinds of pretraining
82 tasks are called Masked Image Modeling (MIM) [5]. There are two main architectures for these
83 models: one that only accesses visible tokens in the encoder and attaches an extra decoder [15, 5],
84 and another that passes both visible and mask tokens into the encoder and has a single linear layer as
85 a decoder [38]. Our formulation is based on the first type of architecture. These mask models serve
86 as a pretrain model, and for downstream tasks, we either finetune or perform linear probing. For
87 classification, a linear head is appended, and the parameters are initialized from the pretrain models.
88 The difference between finetuning and linear probing is that the parameters of the pretrained model
89 are frozen in linear probing.

90 **2.2 Theoretical Analysis of Mask Models**

91 Previous works on mask models have provided theoretical frameworks for understanding the attention
 92 operation in the encoder [6], proposed that MIMs are learning semantics [27], proved a downstream
 93 performance bound for linear probing with MSE loss [22], and claimed that mask models learn global
 94 features that are occlusion invariant [20]. Our work is distinct from these previous works in that
 95 we emphasize the connection between mask modeling and Contrastive Learning. One work also
 96 mentioned that the decoder in MIMs is performing low-rank recovery, but the authors did not link
 97 this to the success of MIMs [6].

98 **2.3 Spectral Contrastive Loss**

99 The Spectral Contrastive loss was proposed as a way to provide a provable guarantee for downstream
 100 task performance [14, 2]. However, some later work has identified issues with the formulation
 101 and stronger assumptions are needed to achieve the guarantee [29]. Despite this, the theoretical
 102 framework that connects Contrastive Learning and Spectral Clustering is still attractive. Our work is
 103 inspired by this analysis framework, but with several differences. In their work, the graph used for
 104 Spectral Clustering is inherent, and the authors argue that matrix factorization approximates the node
 105 representations of the graph. Our work, instead, explicitly writes out the adjacency matrix of a graph
 106 and shows that it is related to the MIM problem. Additionally, our work highlights the importance of
 107 rank, which is often overlooked in previous works.

108 **3 Preliminary and Notations**

109 **3.1 Notations of Masked Autoencoder**

110 Our analysis mainly focus on Masked Autoencoder (MAE) style encoder-decoder structure, where the
 111 input size of encoder is smaller than that of decoder [15]. Denote the encoder in the mask modeling
 112 by f , and the decoder by g , the sampled visible subset by X , and the masked part of the original
 113 image by X_0 . We adopt the Transformer architecture as backbone, where f and g don't change the
 114 shape of inputs [35].

115 **Definition 3.1.** To train the masked autoencoder and achieve the best performance can be interpreted
 116 as solving the minimization problem:

$$\operatorname{argmin}_{f,g,X} \|g \circ f(X) - X_0\|_F^2, \quad (1)$$

117 where $X \in \mathbb{R}^{N \times F}$, $X_0 \in \mathbb{R}^{N_0 \times F}$. We reshape the matrix of image so that N is the number of visible
 118 patches and N_0 is the number of masked patches. We also have the loss defined as

$$\mathcal{L}_{MAE}(f, g, X) = \|g \circ f(X) - X_0\|_F^2 \quad (2)$$

119 **Definition 3.2.** We further define the token mixing layer with weight $W \in \mathbb{R}^{N \times N}$, which mixes
 120 features on the patch level. of inputs. In transformers, the layer is the softmaxed query-key matrix.

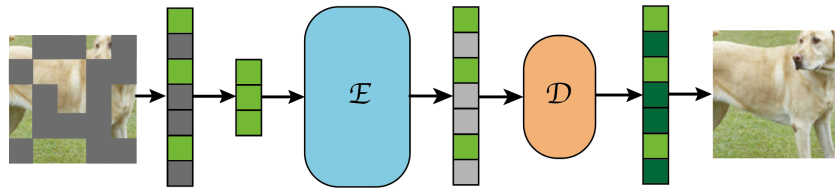


Figure 1: **Overall structure of Mask Autoencoders.** Only visible tokens are passed through the encoder, and a decoder applies encoder features to reconstruct mask tokens.

121 **3.2 Basic Low-rank Recovery Problem**

122 The basic low-rank recovery problem solves the following optimization problem:

$$\operatorname{argmin}_{\hat{D}} \|\hat{D} - D\|_F \text{ subject to } \operatorname{rank}(\hat{D}) \leq \operatorname{rank}(D) \quad (3)$$

123 Based on the Eckart–Young–Mirsky theorem [17], the low-rank approximation problem has a
 124 solution in terms of singular value decomposition of the original matrix, which is in the form:
 125 $\widehat{D} = \sum_{i=1}^r \sigma_i u_i v_i^\top$, where σ_i is the i^{th} singular and u_i and v_i are its corresponding left and right
 126 singular vectors. We could also write it as $\widehat{D} = U_r \Sigma_r V_r^\top$, where U_r, V_r contains the first r columns
 127 of U and V , and Σ_r is an $r \times r$ matrix with the top r leading singular values as diagonal.

128 3.3 Spectral Clustering with Normalized Adjacency Matrix

129 Suppose we have a n -node graph G with the adjacency matrix A :

$$A = \begin{bmatrix} w_{11} & \cdots & w_{1n} \\ \vdots & \ddots & \vdots \\ w_{n1} & \cdots & w_{nn} \end{bmatrix} \quad (4)$$

130 The normalized adjacency matrix is defined as $\mathcal{A} = D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$, where D is the degree matrix of
 131 graph G , which is a diagonal matrix such that $D_{ii} = \sum_{j=1}^n w_{ij} = w_i$. To get k cluster of the graph,
 132 take the leading top k eigenvectors of \mathcal{A} as node embedding and perform k-means algorithm on the
 133 node embedding [26].

134 4 Mask Modeling as a Patch-wise Contrastive Learning

135 4.1 Mask Modeling is Low Rank Recovery

136 In this section, we formalize MAE as a low-rank recovery task. As X is a smaller portion of the
 137 original image and f doesn't change the size of input, $f(X)$ naturally has a lower rank compared to
 138 X_0 . we assume the following condition is true:

139 **Assumption 4.1.** $g \circ f(X)$ has a lower rank than X_0 .

140 *Remark 4.2.* This assumption follows Lemma 6.1 in [6], where they prove an upper bound for the
 141 reconstruction with a low rank decoder assumption. Also in practice, even if g is a non-linear function
 142 that doesn't guarantee low rank assumption, we find that $g(f(X))$ still has a very low rank. Arguably
 143 it is because reconstructing unseen tokens is very hard to optimize and only leading singular vectors
 144 can be approximated.

145 Under this assumption, the minimization problem can be rewritten as

$$\begin{aligned} & \underset{f, g, X}{\operatorname{argmin}} \|g \circ f(X) - X_0\|_F^2 & (5) \\ & \text{subject to } \operatorname{rank}(g \circ f) < \operatorname{rank}(X_0). \end{aligned}$$

146 4.2 Mask Modeling is a Parametric Version of Spectral Clustering

147 In section 3.2, it is showed that the low-rank approximation problem is solved by singular value
 148 decomposition of the higher-ranked matrix. Suppose the required rank is k , then the optimal solution
 149 is a linear combination of top k eigenvectors of $X_0 X_0^\top$ obtained from singular value decomposition.

150 Consider Spectral Clustering which clusters a graph into k connected components such that there
 151 is minimal effect on graph Laplacian. Spectral clustering performs dimensional reduction with k
 152 eigenvectors corresponding with the largest k eigenvalues of the normalized adjacency matrix.

153 Both mask modeling and Spectral Clustering are utilize k eigenvectors, hence, we propose that the
 154 behaviors of mask modeling is similar to the behaviors of Spectral Clustering. Consequently, the
 155 classifier trained based on mask modeling based f and Spectral Clustering based f gives the same
 156 prediction. Formally we have:

157 **Theorem 4.3.** Define weights of adjacency matrix for graph G as $w_{ij} = \langle X_{0r,i}, X_{0r,j} \rangle$, where $X_{0r,i}$
 158 is the low-rank approximation of the representation for i^{th} patch of X_0 . Given the corresponding
 159 normalized adjacency matrix \mathcal{A} , optimizing mask modeling is equivalent to optimize the following
 160 loss on classification downstream tasks.

$$\mathcal{L}_{\text{spec}}(f, g, X) = \|(g \circ f(X))(g \circ f(X))^\top - \mathcal{A}\|_F^2 \quad (6)$$

161 *Proof.* The SVD of X_0 gives $X_0 = U\Sigma V^\top$, then $A = X_{0r}X_{0r}^\top = U_r\Sigma_r^2U_r^\top$.

Since A is symmetric and D is diagonal, \mathcal{A} is symmetric and SVD of \mathcal{A} has the form of $U_{\mathcal{A}}\Sigma_{\mathcal{A}}U_{\mathcal{A}}^\top$. Plug in A gives

$$\begin{aligned}\mathcal{A} &= D^{-\frac{1}{2}}AD^{-\frac{1}{2}} \\ &= D^{-\frac{1}{2}}U_r\Sigma_r^2U_r^\top D^{-\frac{1}{2}} \\ &= \left(D^{-\frac{1}{2}}U_rD^{\frac{1}{2}}\right)\left(D^{-\frac{1}{2}}\Sigma_r^2D^{-\frac{1}{2}}\right)\left(D^{-\frac{1}{2}}U_rD^{\frac{1}{2}}\right)^\top.\end{aligned}$$

162 Therefore, $U_{\mathcal{A}} = D^{-\frac{1}{2}}U_rD^{\frac{1}{2}}$ and $\Sigma_{\mathcal{A}} = D^{-\frac{1}{2}}\Sigma_r^2D^{-\frac{1}{2}}$.

163 The SVD of X_0 can be rewritten as $X_{0r} = D^{\frac{1}{2}}U_{\mathcal{A}}D^{-\frac{1}{2}}\Sigma_rV_r^\top$. With Eckart–Young–Mirsky Theorem, we rewrite the minimization problem of mask modeling as

$$\operatorname{argmin}_{f,g,X} \left\| g \circ f(X) - D^{\frac{1}{2}}U_{\mathcal{A}}D^{-\frac{1}{2}}\Sigma_rV_r^\top \right\|_F^2, \quad (7)$$

165 whose optimal solution is $D^{\frac{1}{2}}U_{\mathcal{A}}D^{-\frac{1}{2}}\Sigma_rV_r^\top$. Note that D is a diagonal matrix, so we could use a decoder to eliminate this term, by setting $g' = D^{-\frac{1}{2}}g$. Therefore, we can discard $D^{\frac{1}{2}}$, making the optimization problem into:

$$\operatorname{argmin}_{f,g,X} \left\| g \circ f(X) - U_{\mathcal{A}}D^{-\frac{1}{2}}\Sigma_rV_r^\top \right\|_F^2. \quad (8)$$

168 With some linear algebra calculation, we could further show that the right hand side of Equation 6 is bounded by the error of Equation 8 with big O notation, i.e. given $\left\| g \circ f(X) - U_{\mathcal{A}}D^{-\frac{1}{2}}\Sigma_rV_r^\top \right\|_F^2 \leq \varepsilon$, $\left\| (g \circ f(X))(g \circ f(X))^\top - \mathcal{A} \right\|_F^2 \leq C\varepsilon$ for some constant C .

171 Let $M = g \circ f(X)$, $N = U_{\mathcal{A}}D^{-\frac{1}{2}}\Sigma_rV_r^\top$, then

$$\begin{aligned}\left\| (g \circ f(X))(g \circ f(X))^\top - \mathcal{A} \right\|_F^2 &= \left\| (g \circ f(X))(g \circ f(X))^\top - (U_{\mathcal{A}}D^{-\frac{1}{2}}\Sigma_rV_r^\top)(U_{\mathcal{A}}D^{-\frac{1}{2}}\Sigma_rV_r^\top)^\top \right\|_F^2 \\ &= \left\| MM^\top - NN^\top \right\|_F^2 \\ &= \frac{1}{4} \left\| (M - N)(M^\top + N^\top) + (M + N)(M^\top - N^\top) \right\|_F^2 \\ &\leq \frac{1}{4} \left(\|M - N\|_F \|M^\top + N^\top\|_F + \|M + N\|_F \|M^\top - N^\top\|_F \right)^2 \\ &= \|M - N\|_F^2 \|M + N\|_F^2\end{aligned}$$

172 Since $N = U_{\mathcal{A}}D^{-\frac{1}{2}}\Sigma_rV_r^\top$ is defined by the original image and M is an approximation to N , we could say that $\|M + N\|_F^2$ is bounded by a constant and thus $\left\| (g \circ f(X))(g \circ f(X))^\top - \mathcal{A} \right\|_F^2 \leq C\varepsilon$, and we finish the proof.

175 □

176 Therefore, MAE learns to approximate the Spectral Clustering features. We further discuss the importance of having appropriate k in Section 5.1

178 4.3 L_{spec} is a Spectral Contrastive Loss

179 Rewrite L_{spec} , mask modeling can be viewed as a token level Contrastive Learning. We define the i^{th} row of $g \circ f(X)$ as $\sqrt{w_i}u_i$, the predicted patch representation. We could rewrite L_{spec} into,

$$\begin{aligned}\mathcal{L}_{spec} &= \left\| (g \circ f(X))(g \circ f(X))^\top - \mathcal{A} \right\|_F^2 \\ &= \left\| (g \circ f(X))(g \circ f(X))^\top - D^{-\frac{1}{2}}AD^{-\frac{1}{2}} \right\|_F^2 \\ &= \sum_{i,j} \left(\frac{w_{ij}}{\sqrt{w_i w_j}} - (\sqrt{w_i}u_i)^\top (\sqrt{w_j}u_j) \right)^2 \\ &= \sum_{i,j} \left(\frac{w_{ij}^2}{w_i w_j} - 2w_{ij}u_i^\top u_j + w_i w_j \cdot (u_i^\top u_j)^2 \right)\end{aligned} \quad (9)$$

181 Apply the kernel trick, changing w_{ij} into W_{ij} , such that $W_{ij} = \exp(\frac{w_{ij}}{2\sigma^2})$ [16]. With a choice of σ ,
 182 we have W_{ij} defined as (or approximates) the probability of u_i and u_j to be a positive pair. Following
 183 the notations of Spectral Contrastive Loss, we make Equation (9) into the form of a Contrastive loss
 184 [14].

$$\mathcal{L}_{\text{spec}} = \mathcal{L}_{\text{cont}} + \text{const}, \quad (10)$$

185 where $\mathcal{L}_{\text{cont}} = -2 \cdot \mathbb{E}_{u, u^+} [u^\top u^+] + \mathbb{E}_{u, u^-} [(u^\top u^-)^2]$

186 *Remark 4.4.* In Haochen et al’s work, u^+ and u^- is defined as positive/negative samples, that has
 187 higher/lower probability to be found in u ’s augmentation space. In mask models, we similarly define
 188 u^+ as patches having higher similarity to u and u^- having lower similarity to u .

189 The above shows that MAE loss is equivalent to a Contrastive loss on masked tokens. We further
 190 show that it inherently perform Contrastive Learning on visible tokens.

191 **Assumption 4.5.** Denote one of the original patch of the i^{th} masked token as $X_{0,i}$, the predicted
 192 feature u_i is a linear combination of features of visible tokens, such that $u_i = \sum_j a_j u'_j$. Assume this
 193 transformation is made by decoder g .

194 **Lemma 4.6.** *Optimal a_j is proportional to patch similarity $u'_j{}^\top X_{0,i}$.*

195 *Proof.* Consider MAE loss with regularization, we have the optimization problem to reconstruct one
 196 patch:

$$\operatorname{argmin}_{\mathbf{a}_j} \left\| \sum_j a_j u'_j - X_{0,i} \right\|_F^2 + \lambda \sum_j a_j^2 \quad (11)$$

197 where $j = 1, 2, \dots, N$, denoting the visible patch representations, and λ is the regularization strength.

198 Let $\mathbf{A} = [a_1, a_2, \dots, a_N]^\top$, $U'_{:,k} = [u'_{1,k}, u'_{2,k}, \dots, u'_{N,k}]^\top$ a rank-1 vector with all patches and k^{th}
 199 latent feature, and $X_{:,k} = [X_{1,k}, X_{2,k}, \dots, X_{N,k}]^\top$ similarly defined. $X_{0,i,k}$ is a scalar corresponds
 200 to the i^{th} patch and k^{th} latent feature. With Assumption 4.5, we decompose it into k rank-1
 201 components.

202 Then Equation 12 becomes:

$$\operatorname{argmin}_{\mathbf{A}} \sum_{k=1}^N \left\| \mathbf{A}^\top U'_{:,k} - X_{0,i,k} \right\|_2^2 + \sum_{k=1}^N \frac{\lambda}{N} \|\mathbf{A}\|_2^2 \quad (12)$$

203 Apply Sherman-Morrison formula [1, 31], which states

$$\left(X + mn^\top \right)^{-1} = X^{-1} - \frac{X^{-1} m n^\top X^{-1}}{1 + n^\top X^{-1} m} \quad (13)$$

204 for any invertible matrix X and rank-1 matrix m and n .

205 Then for each k , the optimal solution for \mathbf{A} would be:

$$\begin{aligned} \hat{\mathbf{A}}_k &= \left(\frac{\lambda}{N} I + U'_{:,k} U'_{:,k}{}^\top \right)^{-1} U'_{:,k} X_{0,i,k} \\ &= N \left(\frac{I}{\lambda} - \frac{\frac{1}{\lambda} U'_{:,k} U'_{:,k}{}^\top \frac{1}{\lambda}}{\frac{1}{N} + U'_{:,k}{}^\top \frac{1}{\lambda} U'_{:,k}} \right) U'_{:,k} X_{0,i,k} \\ &= \frac{N}{\lambda} U'_{:,k} X_{0,i,k} - \frac{N}{\lambda} \frac{U'_{:,k} \left\| U'_{:,k} \right\|_2^2 X_{0,i,k}}{\frac{N}{\lambda} + \left\| U'_{:,k} \right\|_2^2} \end{aligned} \quad (14)$$

206 Since U' is an output of a Layer Norm layer, $\left\| U'_{:,k} \right\|_2^2$ is a constant, denote as C . Plug in Equation 14,
 207 we get the optimal solution for certain k :

$$\hat{\mathbf{A}}_k = \frac{N^2}{\lambda N + \lambda C} U'_{:,k} X_{0,i,k} \quad (15)$$

208 Since Equation 12 solves a sum of least square, the optimal \hat{A} is the mean of \hat{A}_k s, i.e. $A =$
 209 $\frac{N}{\lambda N + \lambda C} U'^T X_{0,i}$, and each optimal a_j is given by

$$\hat{a}_j = \frac{N}{\lambda N + \lambda C} u_j'^T X_{0,i} \quad (16)$$

210 □

211 Meanwhile, as there's only one MLP layer between u'_j and X_j , we view u'_j as an approximation to
 212 X_j , thus

$$\hat{a}_j \approx \frac{N}{\lambda N + \lambda C} X_j X_{0,i} \quad (17)$$

213 which is proportional to the non-parametric patch similarity defined by the original image. Here we
 214 see the representation of masked tokens is mainly composed of similar tokens, performing Contrastive
 215 Learning on masked tokens inherently performs Contrastive Learning on visible tokens.

216 *Remark 4.7.* We could view the layer normalization layer in MAE's decoders as a token level batch
 217 normalization, and the entire decoder as a non-linear projection layer in Contrastive Learning methods
 218 [3].

219 *Remark 4.8.* Though reconstructing masked tokens make it more complicated and less explainable, it
 220 is required as reconstruction visible tokens could lead to a shortcut solution of identity mapping.

221 5 Patch-wise Contrastive Learning Explains Mask Model Behaviors

222 Based on the theoretical framework proposed, we could explain several parameter choice and
 223 architecture design for mask models.

224 5.1 Mask Ratio for Different Modalities

225 In Section 4.2, we demonstrate that mask models are a parametric version of Spectral Clustering, and
 226 they learn to decrease intra-cluster distances while increasing distances between different clusters
 227 through Contrastive Learning. Therefore, an appropriate number of clusters is a crucial factor that
 228 affects the quality of the features learned. When we consider each cluster has a pseudo-class label,
 229 too few or too many classes can both be indistinct when trying to separate the classes. Therefore, we
 230 define the following:

231 **Definition 5.1.** Let s be the ratio of appropriate cluster numbers to total number of tokens. We have

$$s = \frac{\text{num_cluster}}{\text{num_tokens}} = \frac{k}{N_0}, \quad (18)$$

232 where k is the number of leading eigenvectors in Spectral Clustering, and N_0 is the number of tokens
 233 for reconstruction.

234 In mask models, k is subjected to $\text{rank}(g \circ f(X))$, which is determined by the number of visible
 235 tokens N . If we assume $\text{rank}(g \circ f(X))$ is proportional to N , we have:

$$s \propto \frac{N}{N_0} \quad (19)$$

236 As $\frac{N}{N_0} = 1 - \text{mask_ratio}$, s is thus determined by the masking ratio.

237 Intuitively we know that s is smaller for modalities with lower information density, such as video,
 238 vice versa. Therefore, we need a higher masking ratio for lower-density modalities and a smaller one
 239 for higher-information-density modalities [18, 36, 15, 34].

240 5.2 Linear Probing Mask Image Models

241 When tuning MIMs on image classification tasks, there is a significant gap between linear probing
 242 and finetuning [15]. A trick that is often used to improve linear probing performance is to append a
 243 batch normalization layer before the linear head [9]. Without the BN layer, and with an appropriate
 244 batch size, the classification accuracy can drop significantly [37].

245 We argue that this is due to the nature of token-level Contrastive Learning. MIMs only learn to create
 246 and separate several clusters, but do not learn which cluster is indicative of the class label. It is often
 247 the case that the class token from a pretrained MIM does not learn the correct cluster. Therefore,
 248 partially finetuning a token mixing layer can greatly boost accuracy [15]. We also argue that the BN
 249 layer adds non-linearity that partly serves as a token mixing layer. Therefore, we may need to rethink
 250 whether linear probing is a "fair" method to evaluate MIMs.

251 5.3 Network Architecture Matters

252 As discussed in Section 5.2, MIMs do not know how to select important tokens without finetuning.
 253 Therefore, a network architecture with token mixing layers is crucial for the success of mask models
 254 on classification tasks. Finetuning these layers allows MIMs to understand what are important
 255 tokens. Another important architecture design is the number of attention heads. More attention heads
 256 generally improves the expressive ability of mask models by offering more choices of candidate
 257 clusters.

258 6 Experiments

259 We have verified several of our assumptions and mathematical formulations with MAE models
 260 pretrained on Cifar10 and ImageNet-1K (IN-1K) datasets [21, 11]. The model backbones used
 261 for Cifar10 and ImageNet are ViT-Tiny and ViT-Base, respectively [13]. For ViT-Base on IN-1K,
 262 we followed the settings of MAE [15]. The parameters of ViT-Tiny on Cifar-10 are given in the
 263 supplementary materials. We have used the same parameters for linear probing and finetuning, as we
 264 found that changing the optimizer of linear probing to AdamW gives better performance [25].

265 6.1 Low-rank Approximation of Different Mask Ratio

266 To verify Assumption 4.1, and our claim in Section 5.1, we computed the average distance of matrix
 267 factorization components between the reconstructed image and the original image on the Cifar10
 268 dataset. Specifically, given the matrix factorization of the original image and reconstructed image:

$$269 X_0 = \sum_{i=1}^r \sigma_{0,i} u_i v_i^\top \text{ and } \hat{X}_0 = \sum_{i=1}^r \hat{\sigma}_{0,i} \hat{u}_i \hat{v}_i^\top \text{ respectively, we compute } \frac{\|\sigma_{0,i} u_i v_i^\top - \hat{\sigma}_{0,i} \hat{u}_i \hat{v}_i^\top\|_2^2}{\|\sigma_{0,i} u_i v_i^\top\|_2^2}$$

270 for $i = 1, 2, 3, 4, 5$ and models with 50%, 60%, 70%, 80%, 90% mask ratio, shown in Figure 2. Our
 271 results demonstrate: 1. The reconstructed image is low-rank, since the distances are increasing with
 272 index, which means that the leading components are more closely approximated. 2. A higher mask
 273 ratio tends to have smaller number of clusters, as the slope is steeper when mask ratio is higher.

274 6.2 Visualizing cluster results of MAE features

275 In Section 4.3, we argued that MAE is a token level Contrastive Learning, which could result in
 276 good zero-shot segmentation results. Here we demonstrate a few non-cherry-picked segmentation
 277 examples by performing K-means on MAE features in Figure 2.

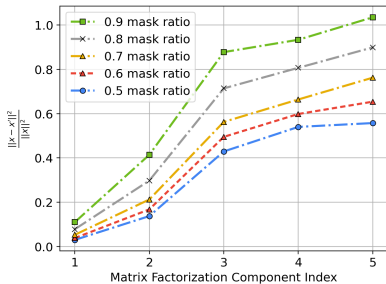


Figure 2: Distances between leading matrix factorization components of reconstructed image and original image.

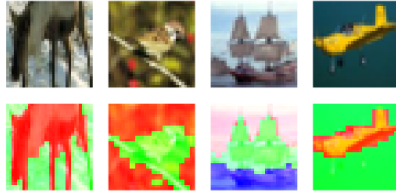


Figure 3: K-means on MAE encoder features.

278 **6.3 Different Probing Methods**

279 In Table 1, we compare three different probing methods mentioned in Section 5.2. We conducted
280 our experiments on both Cifar10 and IN-1K datasets with linear probing with a linear head (LP),
281 non-learnable batch normalization + linear probing (BN + LP), and partial finetuning (Partial FT).
282 For partial finetuning, we tune a linear head and the last qkv projection layer in the encoder, which is
283 also linear.

284 We observe that a non-learnable BN layer can significantly improve the linear probing performance
285 while a learnable token mixing layer can further improve the performance. The performance gain
286 here is larger than what is typically seen in other Contrastive Learning models [9]. This phenomenon
287 provides support for our assumption that the class token learned through the MIM pretext task may
288 not accurately select clusters highly relevant to class labels. The incorporation of token mixture layers
289 enables the reselection of clusters and leads to a performance boost.

Dataset	LP	BN+LP	Partial FT
Cifar10	64.4	76.6 (+12.2)	83.3 (+6.7)
IN-1K	48.0	68.0 (+20.0)	69.3 (+1.3)

Table 1: Classification accuracy with different probing methods.

Heads	LP	FT
3	64.4	89.7
6	67.5 (+3.1)	90.2 (+0.5)

Table 2: Classification accuracy with different attention head numbers.

290 **6.4 Ablation on Number of Attention Heads**

291 This experiment supports our assertion in Section 5.3 that the number of attention heads can impact
292 the performance of MIMs. We pretrained a ViT-Tiny on Cifar10 with 3 and 6 attention heads
293 respectively and their downstream task accuracy with linear probing (without BN) and finetuning
294 are shown in Table 2. By increasing the number of heads, we observed a significant improvement in
295 accuracy for both linear probing and finetuning on the classification task.

296 **7 Discussion**

297 While mask modeling can be seen as a variant of Contrastive Learning at the token level, it is much
298 different from traditional Contrastive methods. The primary distinction lies in the definition of
299 positive and negative pairs: mask modeling methods derive these pairs from natural signals, whereas
300 Contrastive methods use externally sourced human knowledge. This fundamental difference also
301 impacts the data augmentation approach used in each method, with mask modeling employing
302 weak augmentation, and Contrastive methods depending on strong augmentations. It thus raises
303 the intriguing question of whether there could be a unified approach to the augmentation process.
304 Furthermore, there’s a compelling need to explore how mask models can be successfully employed in
305 the arena of multimodal self-supervised learning, especially given the challenge that natural signals
306 often do not align across different modalities.

307 **8 Conclusion**

308 In this paper, we propose a theoretical framework for analyzing mask models. We discover that mask
309 modeling is a form of Contrastive Learning at the token level, which may account for its success. Our
310 framework also addresses important questions regarding the behavior of mask models. We hope that
311 our study will offer valuable insights into designing self-supervised learning algorithms and model
312 architectures.

313 **Limitation** While our paper puts forth a theoretical framework to analyze mask models and elucidates
314 their relationship with Contrastive Learning, we remain unable to provide a provable guarantee for
315 the downstream performance of self-supervised pretraining models. The persisting challenge to
316 fully comprehend Contrastive Learning - specifically, the inductive bias inherent in neural networks -
317 continues to apply to mask models as well [30].

318 **References**

- 319 [1] Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What
320 learning algorithm is in-context learning? investigations with linear models. *arXiv preprint*
321 *arXiv:2211.15661*, 2022.
- 322 [2] Sanjeev Arora, Hrishikesh Khandeparkar, Mikhail Khodak, Orestis Plevrakis, and Nikunj
323 Saunshi. A theoretical analysis of contrastive unsupervised representation learning. *arXiv*
324 *preprint arXiv:1902.09229*, 2019.
- 325 [3] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint*
326 *arXiv:1607.06450*, 2016.
- 327 [4] Roman Bachmann, David Mizrahi, Andrei Atanov, and Amir Zamir. MultiMAE: Multi-modal
328 multi-task masked autoencoders. 2022.
- 329 [5] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. BEit: BERT pre-training of image
330 transformers. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=p-BhZSz59o4>.
331
- 332 [6] Shuhao Cao, Peng Xu, and David A Clifton. How to understand masked autoencoders. *arXiv*
333 *preprint arXiv:2202.03670*, 2022.
- 334 [7] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya
335 Sutskever. Generative pretraining from pixels. In *International Conference on Machine*
336 *Learning*, pages 1691–1703. PMLR, 2020.
- 337 [8] Xiaokang Chen, Mingyu Ding, Xiaodi Wang, Ying Xin, Shentong Mo, Yunhao Wang, Shumin
338 Han, Ping Luo, Gang Zeng, and Jingdong Wang. Context autoencoder for self-supervised
339 representation learning. *arXiv preprint arXiv:2202.03026*, 2022.
- 340 [9] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised
341 vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer*
342 *Vision*, pages 9640–9649, 2021.
- 343 [10] Mark A Davenport and Justin Romberg. An overview of low-rank matrix recovery from
344 incomplete observations. *IEEE Journal of Selected Topics in Signal Processing*, 10(4):608–622,
345 2016.
- 346 [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-
347 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern*
348 *recognition*, pages 248–255. Ieee, 2009.
- 349 [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of
350 deep bidirectional transformers for language understanding. 2018.
- 351 [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
352 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,
353 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image
354 recognition at scale. In *International Conference on Learning Representations*, 2021.
- 355 [14] Jeff Z HaoChen, Colin Wei, Adrien Gaidon, and Tengyu Ma. Provable guarantees for self-
356 supervised deep learning with spectral contrastive loss. *Advances in Neural Information*
357 *Processing Systems*, 34:5000–5011, 2021.
- 358 [15] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked
359 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on*
360 *Computer Vision and Pattern Recognition*, pages 16000–16009, 2022.
- 361 [16] Thomas Hofmann, Bernhard Schölkopf, and Alexander J Smola. Kernel methods in machine
362 learning. *The annals of statistics*, 36(3):1171–1220, 2008.
- 363 [17] Roger A. Horn and Charles R. Johnson. *Matrix Analysis*. Cambridge University Press, 1985.
364 doi: 10.1017/CBO9780511810817.
- 365 [18] Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian
366 Metze, and Christoph Feichtenhofer. Masked autoencoders that listen. In *NeurIPS*, 2022.
- 367 [19] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training
368 by reducing internal covariate shift. In *International conference on machine learning*, pages
369 448–456. PMLR, 2015.

- 370 [20] Xiangwen Kong and Xiangyu Zhang. Understanding masked image modeling via learning
371 occlusion invariant feature. *arXiv preprint arXiv:2208.04164*, 2022.
- 372 [21] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
373 2009.
- 374 [22] Jason D Lee, Qi Lei, Nikunj Saunshi, and Jiacheng Zhuo. Predicting what you already know
375 helps: Provable self-supervised learning. *Advances in Neural Information Processing Systems*,
376 34:309–323, 2021.
- 377 [23] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike
378 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining
379 approach. *arXiv preprint arXiv:1907.11692*, 2019.
- 380 [24] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining
381 Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision
382 and Pattern Recognition*, pages 11976–11986, 2022.
- 383 [25] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint
384 arXiv:1711.05101*, 2017.
- 385 [26] Andrew Ng, Michael Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm.
386 *Advances in neural information processing systems*, 14, 2001.
- 387 [27] Jiachun Pan, Pan Zhou, and Shuicheng Yan. Towards understanding why mask-reconstruction
388 pretraining helps in downstream tasks. *arXiv preprint arXiv:2206.03826*, 2022.
- 389 [28] Vishal M Patel, Hien Van Nguyen, and René Vidal. Latent space sparse and low-rank subspace
390 clustering. *IEEE Journal of Selected Topics in Signal Processing*, 9(4):691–701, 2015.
- 391 [29] Nikunj Saunshi, Jordan Ash, Surbhi Goel, Dipendra Misra, Cyril Zhang, Sanjeev Arora, Sham
392 Kakade, and Akshay Krishnamurthy. Understanding contrastive learning requires incorpor-
393 ating inductive biases. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepes-
394 vari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference
395 on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages
396 19250–19286. PMLR, 17–23 Jul 2022. URL [https://proceedings.mlr.press/v162/
397 saunshi22a.html](https://proceedings.mlr.press/v162/saunshi22a.html).
- 398 [30] Nikunj Saunshi, Jordan Ash, Surbhi Goel, Dipendra Misra, Cyril Zhang, Sanjeev Arora, Sham
399 Kakade, and Akshay Krishnamurthy. Understanding contrastive learning requires incorporating
400 inductive biases. In *International Conference on Machine Learning*, pages 19250–19286.
401 PMLR, 2022.
- 402 [31] Jack Sherman and Winifred J. Morrison. Adjustment of an inverse matrix corresponding to a
403 change in one element of a given matrix. *The Annals of Mathematical Statistics*, 21(1):124–127,
404 1950. ISSN 00034851. URL <http://www.jstor.org/stable/2236561>.
- 405 [32] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions
406 on pattern analysis and machine intelligence*, 22(8):888–905, 2000.
- 407 [33] Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas
408 Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. Mlp-mixer:
409 An all-mlp architecture for vision. *Advances in Neural Information Processing Systems*, 34:
410 24261–24272, 2021.
- 411 [34] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. VideoMAE: Masked autoencoders are
412 data-efficient learners for self-supervised video pre-training. In *Advances in Neural Information
413 Processing Systems*, 2022.
- 414 [35] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
415 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. 2017.
- 416 [36] Alexander Wettig, Tianyu Gao, Zexuan Zhong, and Danqi Chen. Should you mask 15% in
417 masked language modeling? *arXiv preprint arXiv:2202.08005*, 2022.
- 418 [37] Jiantao Wu and Shentong Mo. Object-wise masked autoencoders for fast pre-training. *arXiv
419 preprint arXiv:2205.14338*, 2022.
- 420 [38] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han
421 Hu. Simsim: A simple framework for masked image modeling. 2021.

422 [39] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong.
423 Image BERT pre-training with online tokenizer. In *International Conference on Learning*
424 *Representations*, 2022. URL <https://openreview.net/forum?id=ydopy-e6Dg>.