

# CLIP-MAP: STRUCTURED MATRIX MAPPING FOR PARAMETER-EFFICIENT CLIP COMPRESSION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Contrastive Language-Image Pre-training (CLIP) has achieved widely applications in various computer vision tasks, e.g., text-to-image generation, Image-Text retrieval and Image captioning. However, CLIP suffers from high memory and computation cost, which prohibits its usage to the resource-limited application scenarios. Existing CLIP compression methods typically reduce the size of pre-trained CLIP weights by selecting their subset as weight inheritance for further retraining via mask optimization or important weight measurement. However, these select-based weight inheritance often compromises the feature presentation ability, especially on the extreme compression. In this paper, we propose a novel *mapping-based* CLIP compression framework, **CLIP-Map**. It leverages learnable matrices to map and combine pretrained weights by *Full-Mapping with Kronecker Factorization*, aiming to preserve as much information from the original weights as possible. To mitigate the optimization challenges introduced by the learnable mapping, we propose *Diagonal Inheritance Initialization* to reduce the distribution shifting problem for efficient and effective mapping learning. Extensive experimental results demonstrate that the proposed CLIP-Map outperforms select-based frameworks across various compression ratios, with particularly significant gains observed under high compression settings.

## 1 INTRODUCTION

Large-scale language-image pre-training models, e.g., CLIP (Radford et al., 2021), achieve outstanding zero-shot transfer capability, which has been widely applied to various computer vision tasks, such as, text to image generation (Rombach et al., 2022) and scene understanding (Gu et al., 2021; Rao et al., 2022; Liu et al., 2023). However, such models are accompanied by large parameters and heavy computation costs, which restrict their real-world applications and deployments.

Pruning (Frankle & Carbin, 2018) and knowledge distillation (Hinton et al., 2015; Sanh et al., 2019) are two commonly used techniques for compressing multimodal models. Pruning can be broadly divided into two categories. The first is token pruning (Rao et al., 2021; Bolya et al., 2022; Shi et al., 2023b; Cao et al., 2024), which aims to reduce the computational cost (FLOPs) by selecting and removing less informative tokens during inference. The second is model pruning (Zhou et al., 2021; Sun et al., 2023; Frantar & Alistarh, 2023), which focuses on reducing the number of model parameters and consequently reduces FLOPs. In this paper, our discussion focuses on the model pruning methods. Several previous studies (Shi et al., 2023a; Wu et al., 2023; Lin et al., 2024) have extensively explored how to effectively compress CLIP-like models using a combination of pruning and knowledge distillation. Most of these approaches adopt a pruning-retraining pipeline that first selects and prunes unimportant parameters and then applies retraining to recover performance. The key difference between these methods lies in the different important weight measurement methods to select, with various strategies (Shi et al., 2023a; Lin et al., 2024; Wu et al., 2023) proposed to assess the importance of each parameter.

Regardless of the specific pruning strategy employed, pruning is a *select-based* method and inevitably leads to information loss from the pretrained model. Even with the retraining phase, it remains challenging to recover the information loss caused by dropping unimportant parameters or tokens.

Weight initialization plays a critical role in the training process of neural networks. A well-designed initialization scheme (Glorot & Bengio, 2010; He et al., 2015) can lead to more stable training process

054 and improved final performance. Recent studies (Chen et al., 2015; 2021; Wang et al., 2023a;c;  
 055 Xia et al., 2024) have begun to explore initialization schemes with model growth paradigms, where  
 056 a larger model is initialized by inheriting and expanding a pretrained smaller one, aiming to fully  
 057 transfer the knowledge of smaller pretrained model to a larger one, bring a better initialization and  
 058 accelerate the training process of large models. Among these, LiGO (Wang et al., 2023a) represents a  
 059 prominent line of *mapping-based* work that reformulates model growth as an optimization problem  
 060 by introducing learnable mapping parameters. These mapping parameters are optimized to search the  
 061 most effective expansion strategy, enabling efficient transfer of knowledge from the small to the large  
 062 model. LeTs (Xia et al., 2024) extends the work of LiGO by introducing a compact learnable module  
 063 called **LearnGene** (Wang et al., 2022a; 2023b), along with a set of learnable transformations. Given  
 064 a desired target model size, the corresponding transformation is applied to the LearnGene module to  
 065 produce initialization weights for a model of arbitrary size. This design enables scalable and adaptive  
 066 initialization across a wide range of model capacities.

067 Inspired by these works, we aim to combine these works by replacing traditional select-based  
 068 compression strategies with a learnable mapping-based paradigm, construct a **mapping-retraining**  
 069 pipeline to acquire a compact compressed model with fewer information loss. However, directly  
 070 transferring learnable mapping techniques from model expansion to model compression presents  
 071 several challenges. Firstly, almost all model growth techniques use a partial mapping and weight  
 072 inheritance schemes, where a subset of the pretrained smaller model parameters is copied to the  
 073 larger model and learns a mapping for the remaining part of the larger model. However, in model  
 074 compression scenarios, the target parameter matrix is typically smaller than the original one, making  
 075 such a partial mapping and weight inheritance approach inapplicable. Secondly, mapping a larger  
 076 model into a smaller one requires a substantial number of parameters, which introduces significant  
 077 overhead in both storage and computation, while also increasing the complexity of the optimization  
 078 space. Thirdly, mapping-based methods have been mainly studied in the context of unimodal  
 architectures, leaving the mapping of multimodal models such as CLIP partly unexplored.

079 To solve these limitations, we propose **CLIP-Map**, a mapping-retraining multimodal model com-  
 080 pression pipeline. Our CLIP-Map first acquires a better initialization of compressed model using  
 081 learnable mappings, and then retrain the initialized model using knowledge distillation. Specially,  
 082 we use learnable parameter-efficient matrices  $F^{in}$  and  $F^{out}$  obtained by Kronecker Factorization  
 083 to map large model parameters blocks into smaller counterparts using matrix multiplication. And  
 084 using  $L_{depth}$  to linear combination different layers and get a model with fewer layers. Due to the  
 085 effectiveness of weight inheritance proved by prior works (Chen et al., 2021; Sanh et al., 2019; Wu  
 086 et al., 2023), we initialize  $F^{in}$  and  $F^{out}$  as diagonal matrices. This initialization method ensures that  
 087 part of original pretrained parameters is copied after multiplying with  $F^{in}$  and  $F^{out}$ , enabling easier  
 088 optimization. After initializing in this way, we then optimize them to find the best mapping structure.  
 089 We introduce knowledge distillation in the retraining stage, using the cross-entropy loss between  
 090 the logits of the student model and those of the teacher model as soft labels, enabling the student to  
 091 mimic the behavior of the teacher and effectively inherit its knowledge.

092 We apply our mapping-based method to models of varying scales and find superior performance  
 093 compared to select-based method such as TinyCLIP of similar size across multiple benchmarks.  
 094 Notably, our method maintains competitive performance even under extremely high compression  
 095 ratios. Moreover, our approach requires fewer training epochs, highlighting its efficiency in both  
 096 performance and training cost.

097 In summary, our contributions include:

- 098
- 099 1. We propose a mapping-based compression method that, unlike conventional select-based pruning  
 100 approaches, avoids hard parameter removal and better preserves the full information contained  
 101 in the pretrained model and acquires a better-initialized compact model.
- 102
- 103 2. We replace the the select-based pruning method in the pruning-retraining pipeline with our  
 104 mapping-based compression method, constructing a simplified mapping-retraining pipeline with  
 105 less engineering complexity and better performance.
- 106
- 107 3. Our method maintains strong performance even under high compression ratios with fewer  
 training epochs, demonstrating its effectiveness and efficiency in extreme compression scenarios.

108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161

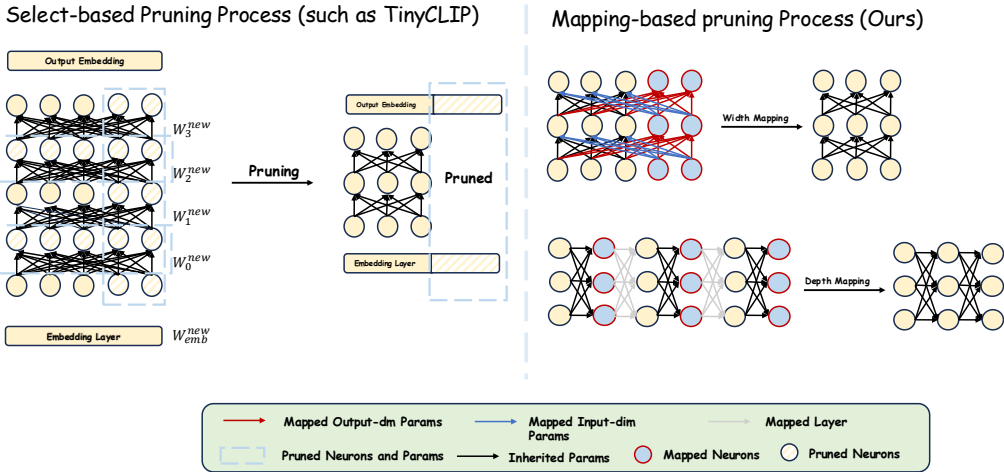


Figure 1: Select-based Compression Method and Mapping-based Growth Method.

## 2 RELATED WORK

### 2.1 MODEL COMPRESSION AND ACCELERATION FOR CLIP

There are many techniques exploring model compression and acceleration, including quantization Krishnamoorthi (2018); Wu et al. (2020), pruning Frankle & Carbin (2018); Sun et al. (2023) and knowledge distillation Hinton et al. (2015); Jiao et al. (2019); Wu et al. (2022); Touvron et al. (2021), etc. With the surge in multimodal foundation models such as CLIP Radford et al. (2021), several studies Gan et al. (2022); Shi et al. (2023a); Touvron et al. (2021) have begun to investigate how these established compression and acceleration techniques on unimodal models can be effectively adapted to the multimodal setting. Unlike unimodal architectures, multimodal models introduce unique challenges such as cross-modal alignment and modality-specific structures, making the direct application of traditional compression methods nontrivial. As a result, recent efforts focus on designing compression frameworks that jointly preserve both modality-specific capabilities and cross-modal interactions.

In the context of CLIP compression and acceleration, the most commonly adopted techniques are pruning and knowledge distillation. A common approach to pruning involves introducing a learnable or heuristic mask to measure the importance of individual weights or channels within the network. The mask can guide us to select a more compact sub-net of the original model, inheriting most original model’s information. Specifically, Gan et al. (2022) shows that the lottery ticket hypothesis (Frankle & Carbin, 2018), originally proposed for unimodal models, still holds in multimodal settings. Regarding inter-modality dependencies, UPop (Shi et al., 2023a) performs joint pruning across both visual and textual modalities. Fig. 1 left provides a simple illustration of the select-based pruning process. Building on pruning, other works (Wang et al., 2022b; Wu et al., 2023; Lin et al., 2024) incorporate knowledge distillation to mitigate information loss introduced by pruning and leverage a powerful teacher model to enhance student performance. Further, CLIP-KD (Yang et al., 2024) extends the work of TinyCLIP(Wu et al., 2023) by investigating an empirical study of CLIP model distillation and identifies optimal design choices for distilling CLIP models. There is another line of pruning work (Rao et al., 2021; Shi et al., 2023b; Cao et al., 2024) focusing on token pruning by eliminating redundant and less informative tokens during the forward pass to achieve inference acceleration. Despite reducing computational overhead during inference through fewer computing tokens, token pruning does not decrease the model’s parameter numbers and may, in some instances, incur additional parameter overhead due to the inclusion of pruning modules.

### 2.2 GROWING PRETRAINED MODELS FOR BETTER INITIALIZATION AND EFFICIENT TRAINING

To mitigate the cost of training deep neural networks, recent approaches (Chen et al., 2015; Gong et al., 2019; Chen et al., 2021; Wang et al., 2023a) explore model growth techniques, which initializes

a larger model using a pretrained smaller model of the same architecture, aiming to transfer the knowledge of pretrained smaller model to the larger model, enabling more efficient training and faster convergence. Net2Net (Chen et al., 2015) introduces function-preserving initialization, which guarantees that a larger model initialized from a smaller one produces identical outputs for the same inputs, thereby preserving and transferring the learned knowledge smoothly. StackBERT (Gong et al., 2019) found that the outputs of different Transformer layers exhibit similarity to some extent. Based on this observation, it proposes to grow the model by simply duplicating a subset of layers from a pretrained BERT (Devlin et al., 2019) and stacking them to form a deeper model. LiGO (Wang et al., 2023a) and LeTs (Xia et al., 2024) formulate model growth as a learnable optimization problem, where a smaller model is transformed into a larger one through parameterized mapping functions as in Fig. 1 right. These mappings are jointly optimized to preserve knowledge while adapting to the increased model capacity.

Our approach is inspired, but distinguished from previous work in the following three major aspects:

- **Mapping-based Compression.** Unlike mapping-based model growth methods, our work focus on model compression using this mapping method.
- **Multimodal Adaptation.** In contrast to previous mapping-based growth methods that are primarily applied to unimodal architectures such as BERT, our approach is specifically designed for multimodal vision-language models like CLIP.
- **Unified and Simplified Pipeline.** Unlike methods that decouple width and depth compression into multiple stages or rely on handcrafted pruning strategies, our framework introduces a unified, end-to-end optimization pipeline. This design simultaneously learns the width and depth compression mappings in a fully differentiable manner, reducing engineering complexity while achieving superior compression-performance trade-offs.

### 3 METHOD

#### 3.1 PRELIMINARIES AND NOTATIONS

**Notation.** Consider a neural network with  $L_1$  layers, each with a hidden dimension of  $D_1$ . Let the weight matrix at layer  $l$  be denoted as  $\mathbf{W}_l \in \mathbb{R}^{D_1 \times D_1}$ . We use an operator  $\mathbf{R}_l$  to acquire a smaller counterpart using matrix multiplication

$$\text{Vec}(\mathbf{W}_l') = \mathbf{R}_l \text{Vec}(\mathbf{W}_l) \in \mathbb{R}^{D_2 \times D_2}, \mathbf{R}_l \in \mathbb{R}^{D_2^2 \times D_1^2}. \quad (1)$$

Here,  $\text{Vec}$  is an operator flattening a parameter matrix into a column vector. We combine each layer’s  $\mathbf{R}_l$  to get the width-compression operator  $\mathbf{R}_{width}$ .

For depth compression, we use a depth-compression operator  $\mathbf{L}_{depth}$  to acquire a shallower network with  $L_2$  layers, where each new layer can be expressed as a linear combination of old layers as in Eq. 2.

$$\mathbf{W}_{l'}^{\text{new}} = \sum_{l=1}^{L_1} \mathbf{L}_{depth}[l', l] \cdot \mathbf{W}_l, \quad l' = 1, \dots, L_2; \mathbf{L}_{depth} \in \mathbb{R}^{L_2 \times L_1}. \quad (2)$$

#### 3.2 THE PROPOSED CLIP-MAP

##### 3.2.1 OVERVIEW

To compress the CLIP model, we need to introduce the operators  $\mathbf{R}_{width}$  and  $\mathbf{L}_{depth}$  operators for the text and visual encoder respectively. Since both modalities share a similar transformer (Vaswani et al., 2017) structure, we focus our discussion on the text encoder for notational simplicity.

We display our mapping-retraining pipeline in Fig. 2. In mapping stage, we freeze original large CLIP model and train the mapping parameters for both image encoder and text encoder. After mapping stage, we acquire a compressed CLIP model inheriting part of the original model. In the second retraining stage, we use this model as initialization of the student model and reuse the original model as a teacher model to distill the student model.

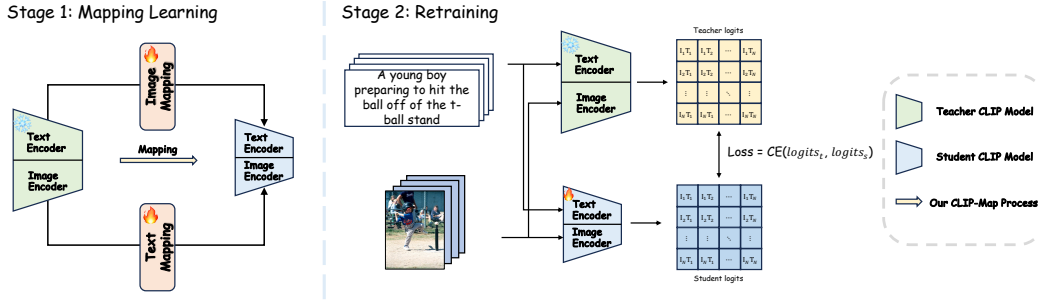


Figure 2: In the mapping-learning stage, we freeze original model’s parameters and train mapping parameters only. In the retraining stage, we use knowledge distillation to distill the student model initialized by mapping stage.

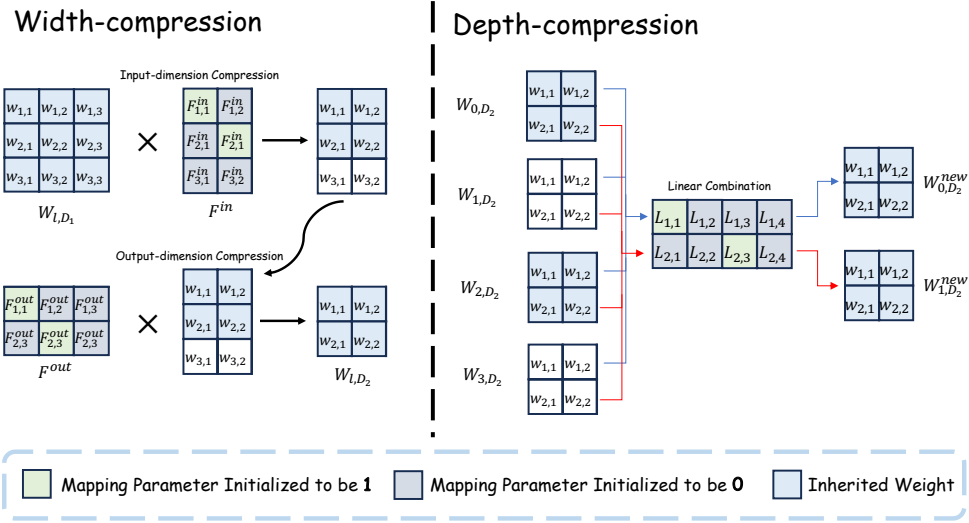


Figure 3: We firstly perform width-compression in both input-dimension and output-dimension on each layers parameter blocks. Then, we perform depth-compression to linear combining the compressed parameter blocks to a new layer parameter block.

In the following sections, we first introduce the core components used in the mapping stage: **Full-Mapping with Kronecker Factorization** and **Diagonal Inheritance Initialization**, which are detailed in Sec 3.2.2 and 3.2.3, respectively. Then, in Sec 3.2.4, we describe the design and training strategy of the retraining stage.

### 3.2.2 FULL-MAPPING WITH KRONECKER FACTORIZATION

As shown in Fig. 3, we apply the Full-Mapping with Kronecker Factorization strategy to compress a pretrained CLIP model. For the parameter block in  $l$ -th layer  $\mathbf{W}_{l,D_1} \in \mathbb{R}^{D_1 \times D_1}$ , our target size is  $\mathbf{W}_{l,D_2} \in \mathbb{R}^{D_2 \times D_2}$ , where  $D_2 < D_1$ . To acquire  $\mathbf{W}_{l,D_2}$ , the simplest approach is using  $\mathbf{R}_l$  and mapping  $\mathbf{W}_{l,D_1}$  using Eq. 1, which leads to the number of parameters  $\mathcal{O}(D_1^2 D_2^2)$ . By leveraging the property of Kronecker products,  $\mathbf{R}_l \text{Vec}(\mathbf{W}_{l,D_1})$  can be reformulated as:

$$\text{Vec}(\mathbf{W}_{l,D_2}) = (\mathbf{F}_l^{\text{in}} \otimes \mathbf{F}_l^{\text{out}}) \text{Vec}(\mathbf{W}_{l,D_1}), \quad (3)$$

$$(\mathbf{F}_l^{\text{in}} \otimes \mathbf{F}_l^{\text{out}}) \text{Vec}(\mathbf{W}_{l,D_1}) = \text{Vec}(\mathbf{F}_l^{\text{out}} \mathbf{W}_{l,D_1} \mathbf{F}_l^{\text{in}^\top}), \quad (4)$$

where  $\otimes$  denotes the Kronecker product and  $\mathbf{F}_l^{\text{in}}, \mathbf{F}_l^{\text{out}} \in \mathbb{R}^{D_2 \times D_1}$ . Therefore, we only use two trainable parameter matrices,  $\mathbf{F}_l^{\text{in}}$  and  $\mathbf{F}_l^{\text{out}}$ , reducing the parameter scale from  $\mathcal{O}(D_1^2 D_2^2)$  to  $\mathcal{O}(D_1 D_2)$ .

Due to the property of the Kronecker product of Eq. 4, it is unnecessary to explicitly construct the full mapping matrix. Instead, we treat the Kronecker factors  $\mathbf{F}_l^{in}$  and  $\mathbf{F}_l^{out}$  as transformations along the input and output dimensions, respectively, and apply them sequentially through standard matrix multiplication.

### 3.2.3 DIAGONAL INHERITANCE INITIALIZATION

During training of mapping, we observe that it is difficult to optimize the mapping matrices  $\mathbf{F}_l^{in}$  and  $\mathbf{F}_l^{out}$ . Since  $\mathbf{F}_l^{in}$  and  $\mathbf{F}_l^{out}$  are derived from a Kronecker decomposition of the original mapping matrix  $\mathbf{R}_{width}$ , directly initializing them using common strategies (e.g. Xavier (Glorot & Bengio, 2010) and Kaiming (He et al., 2015) initialization) often results in *distribution shifting* problem, leading to numerical instability and poor convergence behavior.

**Formulation.** Consider  $\mathbf{A} \in \mathbb{R}^{D_2 \times D_1}$  and  $\mathbf{B} \in \mathbb{R}^{D_2 \times D_1}$  being two matrices independently initialized. Their Kronecker product is denoted as  $\mathbf{R} = \mathbf{A} \otimes \mathbf{B} \in \mathbb{R}^{D_2^2 \times D_1^2}$ . For any element of  $\mathbf{R}$ , we have:

$$\mathbf{R}_{(i-1)D_2+k, (j-1)D_1+l} = \mathbf{A}_{ij} \cdot \mathbf{B}_{kl}, \quad (5)$$

where  $1 \leq i, k \leq D_2$  and  $1 \leq j, l \leq D_1$ .

Assuming each element of  $\mathbf{A}$  and  $\mathbf{B}$  is sampled independently from a zero-mean distribution with variance  $\sigma_A^2$  and  $\sigma_B^2$  respectively, we obtain:

$$\mathbb{E}[\mathbf{R}_{ij}] = \mathbb{E}[\mathbf{A}_{ij} \cdot \mathbf{B}_{kl}] = \mathbb{E}[\mathbf{A}_{ij}] \mathbb{E}[\mathbf{B}_{kl}] = 0, \quad (6)$$

$$\text{Var}(\mathbf{R}_{ij}) = \text{Var}(\mathbf{A}_{ij} \cdot \mathbf{B}_{kl}) = \mathbb{E}[\mathbf{A}_{ij}^2] \mathbb{E}[\mathbf{B}_{kl}^2] = \sigma_A^2 \cdot \sigma_B^2. \quad (7)$$

This implies that although the Kronecker-structured mapping  $\mathbf{R}$  remains zero-mean under independent initialization, its variance becomes multiplicative in nature. If  $\mathbf{A}$  and  $\mathbf{B}$  are both initialized using schemes with moderate variance, the resulting transformation  $\mathbf{R}$  may have a significantly shifting variance:

$$\text{Var}(\mathbf{R}) = \sigma_A^2 \cdot \sigma_B^2. \quad (8)$$

Such uncontrolled variance scaling can lead to unstable optimization dynamics, including gradient vanishing or explosion during the early stage of training. Therefore, careful design or reparameterization of the initialization process is required to ensure stable convergence.

Weight inheritance has been proved to be an effective initialization strategy that facilitates knowledge transfer and accelerates convergence during training (Chen et al., 2021; Sanh et al., 2019; Wu et al., 2023). Thus, we propose **Diagonal Inheritance Initialization**. As illustrated in Eq. 4, the Kronecker factorization can be interpreted as a compression scheme that operates along both the input (in-dim) and output (out-dim) dimensions. Therefore, if we initialize the diagonal elements of both  $\mathbf{F}_l^{in}$  and  $\mathbf{F}_l^{out}$  to 1, it effectively preserves part of the original parameter structure as shown in Fig. 3, thus implementing weight inheritance with much less distribution shifting. Specially, we initialize the diagonal entries of both  $\mathbf{F}_l^{in}$  and  $\mathbf{F}_l^{out}$  and set the off-diagonal elements to zero or small random values

$$(\mathbf{F}_l^{in})_{ij}, (\mathbf{F}_l^{out})_{ij} = \begin{cases} 1, & \text{if } i = j, \\ 0, & \text{otherwise} \end{cases} \quad (9)$$

This structured initialization ensures that the initial Kronecker product approximates an identity-like transformation, i.e.,

$$\mathbf{R}_{width} \approx \mathbf{I}, \quad (10)$$

thus preserving the original parameter semantics in early training. The initialization strategy for  $\mathbf{F}_{in}$  and  $\mathbf{F}_{out}$  across different components of the network is illustrated in detail in A.3.

### 3.2.4 RETRAINING STAGE

In the retraining stage, we introduce knowledge distillation. A CLIP model first extract image embedding and text embedding of a image-text pair using image encoder and text encoder, and then compute image-text similarity(logits). We can use teacher logits as soft labels, forcing student model to mimic the distribution of teacher logits using cross-entropy loss (Wu et al., 2023; Yang et al., 2024).

$$\mathcal{L}_{distill} = CE(\text{logits}_{I2T}^s, \text{logits}_{I2T}^t) + CE(\text{logits}_{T2I}^s, \text{logits}_{T2I}^t). \quad (11)$$

Table 1: Zero-shot image-text retrieval results on MSCOCO and Flickr30K datasets (TR@K: image-to-text retrieval, IR@K: text-to-image retrieval). Our method achieves strong performance with fewer parameters. “2x25ep/3x25ep” in † denotes the method employs the two-stage/three-stage progressive compression strategy, each stage containing 25 training epochs. Detailed architectural configurations for all models can be found in A.3 Tab. 6.

Method	Params (M)	Training Dataset	MSCOCO						Flickr30K						
			TR@1	TR@5	TR@10	IR@1	IR@5	IR@10	TR@1	TR@5	TR@10	IR@1	IR@5	IR@10	
SLIP (Mu et al., 2022)	86+38	YFCC-15M	31.1	-	-	20.3	-	-	57.6	-	-	40.1	-	-	-
CLIP (Wu et al., 2023)	86+38	YFCC-15M	26.5	-	-	17.1	-	-	51.6	-	-	32.2	-	-	-
CLIP (Radford et al., 2021)	86+38	WIT-400M	49.6	74.0	82.1	30.1	54.3	65.1	79.2	95.0	98.0	59.2	83.2	89.6	
OpenCLIP (Cherti et al., 2023; Ilharco et al., 2021)	86+38	LAION-2B	60.0	82.6	89.1	41.3	66.5	76.2	86.2	98.0	99.5	69.8	90.4	94.6	
MetaCLIP (Xu et al., 2023)	86+38	-	58.6	81.0	88.5	41.0	66.7	76.6	-	-	-	-	-	-	
CLIP (Radford et al., 2021)	38+38	WIT-400M	49.3	-	-	30.7	-	-	-	-	-	59.2	-	-	
CLIP (Radford et al., 2021)	86+38	WIT-400M	50.8	-	-	28.3	-	-	81.2	-	-	58.2	-	-	
<i>Training data: YFCC15M</i>															
<b>Compression Ratio: 1.0%</b>															
TinyCLIP (Wu et al., 2023)	0.8+0.3	YFCC15M	10.5	26.1	36.2	5.9	16.9	25.1	21.3	43.3	55.6	11.8	30.2	40.4	
† TinyCLIP (Wu et al., 2023) (3x25ep)	0.8+0.3	YFCC15M	12.5	29.3	39.0	6.9	19.2	28.1	24.5	48.7	60.6	14.6	34.7	46.5	
<b>CLIP-Map<sub>Dist</sub> (Ours)</b>	0.8+0.3	YFCC15M	<b>15.8</b>	<b>34.7</b>	<b>45.3</b>	<b>8.2</b>	<b>22.3</b>	<b>31.7</b>	<b>30.3</b>	<b>56.6</b>	<b>66.3</b>	<b>17.9</b>	<b>40.9</b>	<b>52.6</b>	
<b>Compression Ratio: 10.0%</b>															
TinyCLIP (Wu et al., 2023)	8+3	YFCC15M	33.8	59.6	70.4	20.2	42.9	54.6	60.3	84.0	89.8	40.2	69.0	78.7	
† TinyCLIP (Wu et al., 2023) (2x25ep)	8+3	YFCC15M	36.2	62.6	72.4	21.5	45.4	57.6	62.3	86.1	92.2	42.3	71.5	80.9	
<b>CLIP-Map<sub>Small</sub> (Ours)</b>	8+3	YFCC15M	<b>38.4</b>	<b>64.6</b>	<b>74.4</b>	<b>24.3</b>	<b>48.4</b>	<b>59.9</b>	<b>66.0</b>	<b>88.5</b>	<b>92.6</b>	<b>46.9</b>	<b>73.9</b>	<b>82.6</b>	
<b>CLIP-Map<sub>Small</sub> (Ours, Meta-CLIP)</b>	8+3	YFCC15M	34.3	60.1	70.8	22.6	46.5	58.5	-	-	-	-	-	-	
<b>Compression Ratio: 50.0%</b>															
TinyCLIP (Wu et al., 2023)	39+19	YFCC15M	54.9	79.4	87.2	38.9	64.2	74.1	84.6	96.7	99.0	66.7	88.7	93.7	
<b>CLIP-Map<sub>Base</sub> (Ours)</b>	39+19	YFCC15M	<b>55.1</b>	<b>78.8</b>	<b>86.5</b>	<b>37.9</b>	<b>63.8</b>	<b>74.1</b>	<b>81.9</b>	<b>96.2</b>	<b>98.8</b>	<b>67.6</b>	<b>89.0</b>	<b>94.0</b>	
<b>CLIP-Map<sub>Base</sub> (Ours, Meta-CLIP)</b>	39+19	YFCC15M	53.0	76.9	85.1	37.1	63.4	73.7	-	-	-	-	-	-	
<b>CLIP-Map<sub>Base</sub> (Ours, ResNet-50, wo Retraining)</b>	19+19	YFCC15M	25.5	49.4	61.4	14.3	34.0	45.8	-	-	-	-	-	-	

Table 2: Zero-shot classification top-1 accuracy on 21 downstream datasets. CLIP-Map exhibits strong performance across most tasks.

Method	Image Encoder	Food101	CIFAR10	CIFAR100	SUN397	Stanford Cars	FGVC Aircraft	VOC2007	DTD	Oxford Pets	Caltech101	Flowers102	MNIST	STL10	EuroSAT	RESISC45	GTSRB	KITTI	Country211	PCam	Rendered SST2	ImageNet-1K
<i>Zero-shot performance</i>																						
CLIP (Radford et al., 2021)	ViT-B/16	85.0	88.2	62.0	63.2	57.5	20.0	78.0	40.4	85.2	81.6	65.3	69.5	97.6	43.1	53.5	41.2	31.5	20.2	60.2	60.5	64.4
OpenCLIP (Cherti et al., 2023; Ilharco et al., 2021)	ViT-B/16	86.6	94.9	76.8	70.8	88.5	27.0	78.8	56.3	90.5	83.8	71.3	65.8	97.9	53.4	62.8	48.2	17.0	20.3	56.4	59.9	70.2
† TinyCLIP (Wu et al., 2023)	ViT-0.8M/16	22.2	31.2	11.5	30.1	2.2	3.0	44.6	11.3	17.4	40.5	34.4	11.4	63.8	16.5	12.4	4.3	28.0	4.9	50.0	49.3	16.6
<b>CLIP-Map<sub>Dist</sub> (Ours)</b>	ViT-0.8M/16	<b>24.5</b>	<b>41.9</b>	<b>13.6</b>	<b>32.7</b>	<b>2.1</b>	<b>3.2</b>	<b>38.6</b>	<b>13.7</b>	<b>18.8</b>	<b>41.2</b>	<b>37.7</b>	<b>10.7</b>	<b>72.4</b>	<b>8.4</b>	<b>14.9</b>	<b>3.9</b>	<b>32.5</b>	<b>5.5</b>	<b>50.1</b>	<b>50.0</b>	<b>19.0</b>
† TinyCLIP (Wu et al., 2023)	ViT-8M/16	59.6	72.8	42.1	56.4	7.7	6.9	62.0	28.8	46.2	71.7	58.0	9.8	92.5	23.3	20.7	10.8	15.3	11.7	52.8	50.0	41.1
<b>CLIP-Map<sub>Small</sub> (Ours)</b>	ViT-8M/16	<b>62.8</b>	<b>77.7</b>	<b>45.4</b>	<b>57.6</b>	<b>10.7</b>	<b>9.6</b>	<b>68.5</b>	<b>30.0</b>	<b>50.9</b>	<b>71.9</b>	<b>57.7</b>	<b>9.8</b>	<b>92.3</b>	<b>29.5</b>	<b>27.3</b>	<b>11.0</b>	<b>13.2</b>	<b>12.5</b>	<b>57.4</b>	<b>50.1</b>	<b>42.7</b>
TinyCLIP (Wu et al., 2023)	ViT-39M/16	82.7	91.3	67.7	69.2	51.7	15.1	76.9	47.3	80.8	81.9	70.0	38.0	97.3	52.4	54.9	31.6	11.1	18.3	60.5	49.9	63.5
<b>CLIP-Map<sub>Base</sub> (Ours)</b>	ViT-39M/16	<b>82.7</b>	<b>91.4</b>	<b>68.3</b>	<b>69.2</b>	<b>50.8</b>	<b>22.2</b>	<b>77.0</b>	<b>48.5</b>	<b>83.1</b>	<b>81.4</b>	<b>70.2</b>	<b>13.0</b>	<b>97.3</b>	<b>52.6</b>	<b>55.1</b>	<b>27.0</b>	<b>10.3</b>	<b>18.6</b>	<b>51.7</b>	<b>48.7</b>	<b>63.7</b>

For hard label, we adopt the InfoNCE loss (He et al., 2020; Radford et al., 2021) used in standard CLIP training.

$$\mathcal{L}_{task} = CE(\text{logits}_{I2T}^s, \text{labels}) + CE(\text{logits}_{T2I}^s, \text{labels}). \quad (12)$$

The overall training loss is a weighted sum of hard task loss and distillation soft loss, with the weighting controlled by a coefficient  $\lambda$ .

$$\mathcal{L}_{total} = (1 - \lambda)\mathcal{L}_{task} + \lambda\mathcal{L}_{soft}. \quad (13)$$

## 4 EXPERIMENT

### 4.1 EXPERIMENTAL SETTINGS

**Architecture.** We adopt a Transformer-base (Vaswani et al., 2017) CLIP architecture and design three variants of our model—**CLIP-Map<sub>base</sub>**, **CLIP-Map<sub>small</sub>** and **CLIP-Map<sub>tiny</sub>**—each tailored to different application scenarios and resource constraints. CLIP-Map<sub>base</sub>, CLIP-Map<sub>small</sub> and CLIP-Map<sub>tiny</sub> are mapped and distilled from OpenCLIP-ViT-B/16 (Cherti et al., 2023; Ilharco et al., 2021). We also evaluate our approach on Meta-CLIP (Xu et al., 2023) and CLIP using ResNet (He et al., 2016) as vision encoder to validate the generalization capability of our approach to any CLIP-like architecture. It should be noted that when ResNet serves as the vision encoder, we limit the process to the Mapping stage 5-epochs training, and don’t perform the subsequent Retraining stage.

**Training setup.** We use YFCC-15M (Li et al., 2021), a large-scale public image-text pair datasets for training. The overall training procedure is organized into two stages: mapping stage and retraining stage. During mapping stage, we optimize learnable mapping matrices, and in retraining stage, we use each mapped model’s original model as a teacher model to distill the mapped model. All models are trained on 32 Nvidia H800. Detailed training settings are presented in A.5.

Table 3: Zero-shot IN-val performance comparison of various compressed CLIP models under different datasets.

Methods	Dataset	# of Seen Samples	Params(M) <i>img+txt</i>	Zero-shot IN-val
†TinyCLIP-0.8M/16	YFCC-15M	1.125B	0.8 + 0.3	16.6
CLIP-Map <sub>tiny</sub>	YFCC-15M	0.45B	0.8 + 0.3	19.0
†TinyCLIP-8M/16 (Wu et al., 2023)	YFCC-15M	0.75B	8 + 3	41.1
ViT-T/16 (Yang et al., 2024)	CC3M (Sharma et al., 2018) + CC12M (Changpinyo et al., 2021)	0.48B	5.6 + 21.3	42.6
CLIP-Map <sub>small</sub>	YFCC-15M	0.45B	8+3	42.7
MobileCLIP-S0 (Vasu et al., 2024)	DataCompDR-12M	0.74B	11.4 + 42.4	59.1
MoPE-CLIP <sub>base</sub> (Lin et al., 2024)	YFCC-15M	0.30B	86+42	60.7
TinyCLIP-39M/16 Wu et al. (2023)	YFCC-15M	0.75B	39+19	63.5
CLIP-Map <sub>base</sub>	YFCC-15M	0.30B	39+19	<b>63.7</b>

Table 4: Effect of different initialization steps on downstream performance. Longer initialization improves accuracy across ImageNet-1K classification (IN-1K) and MSCOCO retrieval.

Init. Steps	IN-1K Top-1 (%)	MSCOCO TR@K			MSCOCO IR@K		
	Acc.	TR@1	TR@5	TR@10	IR@1	IR@5	IR@10
Manual Drop (0 epoch)	41.1	33.8	59.6	70.4	20.2	42.9	54.6
0.28(1000steps) + 25 epochs	39.7	35.2	61.1	71.3	21.7	44.6	56.4
1 + 24 epochs	39.6	35.7	61.5	71.6	21.0	44.3	56.4
3 + 22 epochs	41.9	37.6	63.1	<b>73.9</b>	23.0	46.9	58.7
5 + 20 epochs	<b>42.1</b>	<b>38.3</b>	<b>63.6</b>	<b>73.9</b>	<b>23.1</b>	<b>47.1</b>	<b>59.2</b>
7 + 18 epochs	40.8	36.2	61.9	72.3	22.0	45.5	57.6

**Benchmarks and Metrics.** We evaluate our method on two tasks: zero-shot classification and zero-shot retrieval. For the retrieval task, we adopt the MSCOCO (Lin et al., 2014) and Flickr30K (Plummer et al., 2015) test sets, using top-k image-to-text and text-to-image recall rates (TR@K and IR@K) as evaluation metrics. For the classification task, we mainly conduct evaluations on the ImageNet-1K (Deng et al., 2009) test set, and report top-1 accuracy (Acc@1) to measure model performance. Beyond ImageNet-1K, we further evaluate the classification capabilities of our model on a diverse set of downstream zero-shot classification tasks.

## 4.2 MAIN RESULTS

**Zero-shot Image-text Retrieval.** Tab. 1 presents the zero-shot retrieval performance of our CLIP-Map<sub>tiny</sub>, CLIP-Map<sub>small</sub> and CLIP-Map<sub>base</sub> models on the MSCOCO (5k test set) and Flickr30k (1k test set). To enable a fair comparison, we replicate three TinyCLIP models of equivalent size using the official TinyCLIP approach, under both progressive and non-progressive compression settings. For the progressive compression strategy, we start from the original OpenCLIP-ViT-B/16 and progressively compress it to 50.0%, 10.0%, and 1.0% of its original size. Models obtained via this progressive compression are marked with a † symbol in the result table. Under the compression ratio of 10.0% and 1.0%, our CLIP-Map<sub>small</sub> and CLIP-Map<sub>tiny</sub> consistently outperform TinyCLIP of the same size across all recall metrics. And our CLIP-Map<sub>base</sub> achieves competitive performance but with fewer training epochs.

**Zero-shot Classification.** Tab. 2 shows the zero-shot classification performance of our method on multiple classification datasets. Compared to TinyCLIP models of similar sizes, our approach demonstrates competitive performance at the *base* scale, achieving results comparable to the baseline. Notably, under the *small* and *tiny* configurations, our model consistently outperforms the baseline across the majority of classification benchmarks, with substantial performance gains observed on several datasets.

**Comparison with Other VLM Compression Methods.** As shown in Tab. 3, we compare our method with other CLIP compression methods. The effectiveness of different approaches is evaluated via zero-shot classification accuracy on ImageNet-1K, while the overall efficiency is assessed based on the total volume of seen samples during the entire training pipeline and the model size. Compared to MoPE-CLIP (Lin et al., 2024) and CLIP-KD (Yang et al., 2024), our methods can get a better performance even with smaller model size and fewer seen samples, highlighting our method in both effectiveness and efficiency. MobileCLIP (Vasu et al., 2024) leverages an augmented dataset, DataCompDR (Gadre et al., 2023; Vasu et al., 2024), constructed via advanced data augmentation techniques. Compared to commonly used datasets such as YFCC-15M and CC3M, DataCompDR



offers higher data quality, thereby enabling MobileCLIP to achieve superior accuracy even with a significantly smaller image encoder. The practical training speed-up brought by our method over TinyCLIP is visualized and presented in A.6. Additional comparative results with other compression methods can be found in A.4.

### 4.3 ABLATION STUDIES

To demonstrate the effectiveness of our method, we conduct ablation studies for different initialization methods and mapping/retraining duration. We also investigate the effect of training loss in A.8.

**Different Initialization Methods.** To validate the effectiveness of *Diagonal Inheritance Initialization*, we conduct experiments on the OpenCLIP-ViT-B/16 model under a 10.0% compression ratio, comparing the training process and performance of different initialization strategies for the mapping parameter matrices. We adopt Random init, Kaiming init (He et al., 2015) and Xavier init (Glorot & Bengio, 2010) for comparison. The loss curves and final performance of different initialization methods are illustrated in A.7 Fig. 6 and summarized in Tab. 5. The *Diagonal Inheritance Initialization* lead to a easier optimization and faster convergence loss curve, bringing much better performance on test sets.

Table 5: Effect of different initialization methods on MSCOCO Recall@1 and IN-1K classification.

Init. Steps	IN-1K Top-1 (%)	MSCOCO Recall@1
	Acc.	I→T / T→I
Random Init	0.1	0.02 / 0.01
Kaiming Init (He et al., 2015)	4.4	3.5 / 2.0
Xavier Init (Glorot & Bengio, 2010)	4.9	4.1 / 2.3
<b>Diag Init (Ours)</b>	<b>28.9</b>	<b>28.5 / 15.9</b>

**Mapping/Retraining Duration.** To determine the optimal number of training epochs for the mapping stage and retraining stage, we conduct 10% compression experiments on the OpenCLIP-ViT-B/16 model. The results are presented in Tab. 4 and the visualization of the weight distribution is presented in A.7. As the mapping stage expands, the distribution of the mapping matrix gradually evolves from an initial diagonal pattern toward a more uniform structure, indicating that the optimization process is progressively searching for an optimal compression mapping. Moreover, as the mapping stage is extended, the performance of the final compressed model is consistently improved. However, we observe that an excessively long mapping stage may lead to performance degradation and introduce unnecessary computational overhead. Thus, we adopt a moderate number of training epochs(5ep) during the mapping stage, which we found providing a good balance between initialization quality and final performance, for both CLIP-Map<sub>tiny</sub> and CLIP-Map<sub>small</sub>.

## 5 CONCLUSION

In this paper, we propose **CLIP-Map**, a novel *mapping-based* compression framework for CLIP-like multimodal models. CLIP-Map employs learnable transformation matrices to perform layer-wise width compression via matrix multiplication, and leverages learnable linear combinations of parameter blocks across layers to realize depth compression.

To mitigate the overhead of parameter storage and computation introduced by the learnable mappings, we propose a *Full Mapping with Kronecker Factorization* strategy, which significantly reduces the parameter complexity. Furthermore, to address the optimization difficulty, we introduce a *Diagonal Inheritance Initialization* scheme that not only enables weight inheritance from the pretrained model but also improves training stability and convergence, resulting in a well-initialized compressed model.

In addition, we build a complete *mapping-retraining pipeline*, where we incorporate knowledge distillation during retraining to facilitate knowledge transfer from the teacher model. By aligning the output logits of the student and teacher models, our method ensures performance preservation even under aggressive compression settings.

Experimental results on various zero-shot retrieval and zero-shot classification benchmarks demonstrate that CLIP-Map achieves competitive or superior performance compared to existing baselines, while providing significant reductions in model size and training overhead.

## REFERENCES

- 486  
487  
488 Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy  
489 Hoffman. Token merging: Your vit but faster. *arXiv preprint arXiv:2210.09461*, 2022.
- 490  
491 Jianjian Cao, Peng Ye, Shengze Li, Chong Yu, Yansong Tang, Jiwen Lu, and Tao Chen. Madtp:  
492 Multimodal alignment-guided dynamic token pruning for accelerating vision-language transformer.  
493 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp.  
494 15710–15719, 2024.
- 495  
496 Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing  
497 web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the  
IEEE/CVF conference on computer vision and pattern recognition*, pp. 3558–3568, 2021.
- 498  
499 Cheng Chen, Yichun Yin, Lifeng Shang, Xin Jiang, Yujia Qin, Fengyu Wang, Zhi Wang, Xiao Chen,  
500 Zhiyuan Liu, and Qun Liu. bert2bert: Towards reusable pretrained language models. *arXiv preprint  
arXiv:2110.07143*, 2021.
- 501  
502 Tianqi Chen, Ian Goodfellow, and Jonathon Shlens. Net2net: Accelerating learning via knowledge  
503 transfer. *arXiv preprint arXiv:1511.05641*, 2015.
- 504  
505 Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade  
506 Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for  
507 contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer  
Vision and Pattern Recognition*, pp. 2818–2829, 2023.
- 508  
509 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale  
510 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
511 pp. 248–255. Ieee, 2009.
- 512  
513 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep  
514 bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of  
the North American chapter of the association for computational linguistics: human language  
515 technologies, volume 1 (long and short papers)*, pp. 4171–4186, 2019.
- 516  
517 Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural  
518 networks. *arXiv preprint arXiv:1803.03635*, 2018.
- 519  
520 Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in  
521 one-shot. In *International Conference on Machine Learning*, pp. 10323–10337. PMLR, 2023.
- 522  
523 Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen,  
524 Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, et al. Datacomp: In search of the  
525 next generation of multimodal datasets. *Advances in Neural Information Processing Systems*, 36:  
526 27092–27112, 2023.
- 527  
528 Zhe Gan, Yen-Chun Chen, Linjie Li, Tianlong Chen, Yu Cheng, Shuohang Wang, Jingjing Liu, Lijuan  
529 Wang, and Zicheng Liu. Playing lottery tickets with vision and language. In *Proceedings of the  
AAAI Conference on Artificial Intelligence*, volume 36, pp. 652–660, 2022.
- 530  
531 Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural  
532 networks. In *Proceedings of the thirteenth international conference on artificial intelligence and  
533 statistics*, pp. 249–256. JMLR Workshop and Conference Proceedings, 2010.
- 534  
535 Linyuan Gong, Di He, Zhuohan Li, Tao Qin, Liwei Wang, and Tiejun Liu. Efficient training of  
536 bert by progressively stacking. In *International conference on machine learning*, pp. 2337–2346.  
537 PMLR, 2019.
- 538  
539 Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision  
and language knowledge distillation. *arXiv preprint arXiv:2104.13921*, 2021.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing  
human-level performance on imagenet classification. In *Proceedings of the IEEE international  
conference on computer vision*, pp. 1026–1034, 2015.

- 540 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
541 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
542 pp. 770–778, 2016.
- 543 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for  
544 unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on*  
545 *computer vision and pattern recognition*, pp. 9729–9738, 2020.
- 546 Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv*  
547 *preprint arXiv:1503.02531*, 2015.
- 548 Lu Hou, Zhiqi Huang, Lifeng Shang, Xin Jiang, Xiao Chen, and Qun Liu. Dynabert: Dynamic  
549 bert with adaptive width and depth. *Advances in Neural Information Processing Systems*, 33:  
550 9782–9793, 2020.
- 551 Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori,  
552 Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali  
553 Farhadi, and Ludwig Schmidt. Openclip, July 2021. URL [https://doi.org/10.5281/  
554 zenodo.5143773](https://doi.org/10.5281/zenodo.5143773). If you use this software, please cite it as below.
- 555 Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu.  
556 Tinybert: Distilling bert for natural language understanding. *arXiv preprint arXiv:1909.10351*,  
557 2019.
- 558 Raghuraman Krishnamoorthi. Quantizing deep convolutional networks for efficient inference: A  
559 whitepaper. *arXiv preprint arXiv:1806.08342*, 2018.
- 560 Yangguang Li, Feng Liang, Lichen Zhao, Yufeng Cui, Wanli Ouyang, Jing Shao, Fengwei Yu, and  
561 Junjie Yan. Supervision exists everywhere: A data efficient contrastive language-image pre-training  
562 paradigm. *arXiv preprint arXiv:2110.05208*, 2021.
- 563 Haokun Lin, Haoli Bai, Zhili Liu, Lu Hou, Muye Sun, Linqi Song, Ying Wei, and Zhenan Sun. Mope-  
564 clip: Structured pruning for efficient vision-language models with module-wise pruning error  
565 metric. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
566 pp. 27370–27380, 2024.
- 567 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
568 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer vision–  
569 ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings,  
570 part v 13*, pp. 740–755. Springer, 2014.
- 571 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in*  
572 *neural information processing systems*, 36:34892–34916, 2023.
- 573 Norman Mu, Alexander Kirillov, David Wagner, and Saining Xie. Slip: Self-supervision meets  
574 language-image pre-training. In *European conference on computer vision*, pp. 529–544. Springer,  
575 2022.
- 576 Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and  
577 Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer  
578 image-to-sentence models. In *Proceedings of the IEEE international conference on computer*  
579 *vision*, pp. 2641–2649, 2015.
- 580 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
581 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
582 models from natural language supervision. In *International conference on machine learning*, pp.  
583 8748–8763. PmLR, 2021.
- 584 Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. Dynamicvit:  
585 Efficient vision transformers with dynamic token sparsification. *Advances in neural information*  
586 *processing systems*, 34:13937–13949, 2021.

- 594 Yongming Rao, Wenliang Zhao, Guangyi Chen, Yansong Tang, Zheng Zhu, Guan Huang, Jie Zhou,  
595 and Jiwen Lu. Denseclip: Language-guided dense prediction with context-aware prompting.  
596 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp.  
597 18082–18091, 2022.
- 598 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
599 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-  
600 ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 601 Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of  
602 bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- 603 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi  
604 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An  
605 open large-scale dataset for training next generation image-text models. *Advances in neural  
606 information processing systems*, 35:25278–25294, 2022.
- 607 Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned,  
608 hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th  
609 Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp.  
610 2556–2565, 2018.
- 611 Dachuan Shi, Chaofan Tao, Ying Jin, Zhendong Yang, Chun Yuan, and Jiaqi Wang. Upop: Unified and  
612 progressive pruning for compressing vision-language transformers. In *International Conference  
613 on Machine Learning*, pp. 31292–31311. PMLR, 2023a.
- 614 Dachuan Shi, Chaofan Tao, Anyi Rao, Zhendong Yang, Chun Yuan, and Jiaqi Wang. Crossget:  
615 Cross-guided ensemble of tokens for accelerating vision-language transformers. *arXiv preprint  
616 arXiv:2305.17455*, 2023b.
- 617 Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter. A simple and effective pruning approach for  
618 large language models. *arXiv preprint arXiv:2306.11695*, 2023.
- 619 Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé  
620 Jégou. Training data-efficient image transformers & distillation through attention. In *International  
621 conference on machine learning*, pp. 10347–10357. PMLR, 2021.
- 622 Pavan Kumar Anasosalu Vasu, Hadi Pouransari, Fartash Faghri, Raviteja Vemulapalli, and Oncel  
623 Tuzel. Mobileclip: Fast image-text models through multi-modal reinforced training. In *Proceedings  
624 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15963–15974,  
625 2024.
- 626 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
627 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing  
628 systems*, 30, 2017.
- 629 Peihao Wang, Rameswar Panda, Lucas Torroba Hennigen, Philip Greengard, Leonid Karlinsky,  
630 Rogerio Feris, David Daniel Cox, Zhangyang Wang, and Yoon Kim. Learning to grow pretrained  
631 models for efficient transformer training. *arXiv preprint arXiv:2303.00980*, 2023a.
- 632 Qiu-Feng Wang, Xin Geng, Shu-Xia Lin, Shi-Yu Xia, Lei Qi, and Ning Xu. LearnGene: From  
633 open-world to your learning task. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
634 volume 36, pp. 8557–8565, 2022a.
- 635 Qiufeng Wang, Xu Yang, Shuxia Lin, Jing Wang, and Xin Geng. LearnGene: Inheriting condensed  
636 knowledge from the ancestry model to descendant models. *arXiv preprint arXiv:2305.02279*,  
637 2023b.
- 638 Tiannan Wang, Wangchunshu Zhou, Yan Zeng, and Xinsong Zhang. Efficientvlm: Fast and accurate  
639 vision-language models via knowledge distillation and modal-adaptive pruning. *arXiv preprint  
640 arXiv:2210.07795*, 2022b.
- 641 Yite Wang, Jiahao Su, Hanlin Lu, Cong Xie, Tianyi Liu, Jianbo Yuan, Haibin Lin, Ruoyu Sun, and  
642 Hongxia Yang. Lemon: Lossless model expansion. *arXiv preprint arXiv:2310.07999*, 2023c.

648 Haiyan Wu, Yuting Gao, Yinqi Zhang, Shaohui Lin, Yuan Xie, Xing Sun, and Ke Li. Self-supervised  
649 models are good teaching assistants for vision transformers. In *International Conference on*  
650 *Machine Learning*, pp. 24031–24042. PMLR, 2022.

651  
652 Hao Wu, Patrick Judd, Xiaojie Zhang, Mikhail Isaev, and Paulius Micikevicius. Integer quantization  
653 for deep learning inference: Principles and empirical evaluation. *arXiv preprint arXiv:2004.09602*,  
654 2020.

655 Kan Wu, Houwen Peng, Zhenghong Zhou, Bin Xiao, Mengchen Liu, Lu Yuan, Hong Xuan, Michael  
656 Valenzuela, Xi Stephen Chen, Xinggang Wang, et al. Tinyclip: Clip distillation via affinity  
657 mimicking and weight inheritance. In *Proceedings of the IEEE/CVF International Conference on*  
658 *Computer Vision*, pp. 21970–21980, 2023.

659 Shiyu Xia, Yuankun Zu, Xu Yang, and Xin Geng. Initializing variable-sized vision transformers from  
660 learnene with learnable transformation. *Advances in Neural Information Processing Systems*, 37:  
661 43341–43366, 2024.

662  
663 Hu Xu, Saining Xie, Xiaoqing Ellen Tan, Po-Yao Huang, Russell Howes, Vasu Sharma, Shang-Wen  
664 Li, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. Demystifying clip data. *arXiv*  
665 *preprint arXiv:2309.16671*, 2023.

666 Chuanguang Yang, Zhulin An, Libo Huang, Junyu Bi, Xinqiang Yu, Han Yang, Boyu Diao, and  
667 Yongjun Xu. Clip-kd: An empirical study of clip model distillation. In *Proceedings of the*  
668 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15952–15962, 2024.

669  
670 Aojun Zhou, Yukun Ma, Junnan Zhu, Jianbo Liu, Zhijie Zhang, Kun Yuan, Wenxiu Sun, and  
671 Hongsheng Li. Learning n: m fine-grained structured sparse neural networks from scratch. *arXiv*  
672 *preprint arXiv:2102.04010*, 2021.

## 674 A APPENDIX

675  
676 In the appendix, we present the detailed design and experiment results of our method.

### 677 A.1 USE OF LLMs

678  
679 In this paper, we use LLMs to polish writing and to assist with LaTeX typesetting and formatting.

### 680 A.2 ETHICS AND REPRODUCIBILITY STATEMENT

681  
682 We have reviewed the ICLR Code of Ethics and make sure that our research confirms every respect  
683 with the ICLR Code of Ethics. And we list all the concrete settings required in original paper and  
684 Appendix and make sure all the results are reproducible.

### 685 A.3 DETAILED ALGORITHM DESIGNING AND MODEL CONFIGURATION

686  
687 In our setting, both the vision encoder and text encoder of CLIP adopt Transformer-based architectures,  
688 which exhibit similar structural forms. For simplicity of notation and without loss of generality, we  
689 focus our discussion on the text encoder, as the mapping strategy for the vision encoder follows an  
690 analogous formulation.

691  
692 **Embedding Layer.** For embedding layers, we use learnable matrix  $F_{emb}^{out}$  to map and get a smaller  
693 embedding layer with reduced embedding dimension.

694  
695 **Multi-Head Attention Layer.** For the  $\ell$ -th MHA layer, we have the following components:  
696  $W_l^Q, W_l^K, W_l^V, W_l^O$ . Instead of assigning distinct  $F^{in}$  and  $F^{out}$  for each component, we re-  
697 duce the number of parameters by reusing a shared transformation  $F_{emb}^{out}$  across all components.  
698 Specially, we set  $F_{l,Q}^{in}, F_{l,K}^{in}, F_{l,V}^{in}$  and  $F_{l,O}^{out}$  to be  $F_{emb}^{out}$ , and  $F_{l,O}^{in} = F_{l,V}^{out}$ .

699  
700 **Feed-Forward Layer.** For the  $\ell$ -th FFN layer, we have learnable components  $W_l^{fc1}$  and  $W_l^{fc2}$ . We  
701 set  $F_{l,fc1}^{in} = F_{emb}^{out}$  and  $F_{l,fc2}^{out} = F_{emb}^{out}$ .

Table 6: Model configurations for zero-shot image-text retrieval experiments.

Method	ViT	Vision Encoder		Text Encoder		Params (M)	Training Dataset
		Width	Depth	Width	Depth		
SLIP (Mu et al., 2022)	ViT-B/16	768	12	512	12	86+38	YFCC-15M
CLIP (Wu et al., 2023)	ViT-B/16	768	12	512	12	86+38	YFCC-15M
CLIP (Radford et al., 2021)	ViT-B/16	768	12	512	12	86+38	WIT-400M (Radford et al., 2021)
OpenCLIP (Cherti et al., 2023; Ilharco et al., 2021)	ViT-B/16	768	12	512	12	86+38	LAION-2B (Schuhmann et al., 2022)
MetaCLIP (Xu et al., 2023)	ViT-B/16	768	12	512	12	86+38	-
CLIP (Radford et al., 2021)	ResNet-50 (He et al., 2016)	-	-	512	12	38+38	WIT-400M
CLIP (Radford et al., 2021)	ResNet-101 (He et al., 2016)	-	-	512	12	86+38	WIT-400M
<i>Training data: YFCC15M</i>							
<b>Compression Ratio: 1.0%</b>							
TinyCLIP (Wu et al., 2023)	ViT-B/16	128	4	128	2	0.8+0.3	YFCC15M
† TinyCLIP (Wu et al., 2023) (3×25Sep)	ViT-B/16	128	4	128	2	0.8+0.3	YFCC15M
<b>CLIP-Map<sub>tiny</sub> (Ours)</b>	ViT-B/16	512	12	512	6	0.8+0.3	YFCC15M
<b>Compression Ratio: 10.0%</b>							
TinyCLIP (Wu et al., 2023)	ViT-B/16	256	10	256	3	8+3	YFCC15M
† TinyCLIP (Wu et al., 2023) (2×25Sep)	ViT-B/16	256	10	256	3	8+3	YFCC15M
<b>CLIP-Map<sub>small</sub> (Ours)</b>	ViT-B/16	256	10	256	3	8+3	YFCC15M
<b>CLIP-Map<sub>small</sub> (Ours, Meta-CLIP)</b>	ViT-B/16	256	10	256	3	8+3	YFCC15M
<b>Compression Ratio: 50.0%</b>							
TinyCLIP (Wu et al., 2023)	ViT-B/16	512	24	768	6	39+19	YFCC15M
<b>CLIP-Map<sub>base</sub> (Ours)</b>	ViT-B/16	512	12	512	6	39+19	YFCC15M
<b>CLIP-Map<sub>base</sub> (Ours, Meta-CLIP)</b>	ViT-B/16	512	12	512	6	39+19	YFCC15M
<b>CLIP-Map<sub>base</sub> (Ours, ResNet-50, wo Retraining)</b>	ResNet-19M	-	-	512	6	19+19	YFCC15M

**Training Loss.** In the mapping stage, we train the mapping parameters using standard CLIP task loss as in Eq. 12.

**Detailed Model Configuration.** Detailed model architectures are provided in Tab. 6 to complement the results in Tab. 1.

#### A.4 COMPARISON WITH OTHER METHOD

Table 7: Comparison of image–text retrieval performance.

Method	Dataset	Params (img + text)	I2T@1	I2T@5	I2T@10	T2I@1	T2I@5	T2I@10
UPop (Shi et al., 2023a)	MSCOCO	280.2M	56.1	82.4	90.2	41.1	71.0	81.4
EfficientVLM (Wang et al., 2022b)	CC3M	152M + 42M	46.6	71.7	81.3	35.9	61.6	71.8
DynaCLIP (Hou et al., 2020)	CC3M	86M + 42M	51.3	75.5	84.6	35.8	61.8	72.6
CLIP-Map-base (Ours)	YFCC15M	39M + 19M	55.1	78.8	86.5	37.9	63.8	74.1
CLIP-Map-small (Ours)	YFCC15M	8M + 3M	38.4	64.6	74.4	24.3	48.4	59.9

To further validate the generalization capability of our approach, we compare its performance with other compression methods. We have surveyed several other VLM compression methods, including UPop (Shi et al., 2023a) and EfficientVLM (Wang et al., 2022b) as in Tab.7. We find that these methods don’t report experiments at higher compression ratios, and the compressed models remain relatively large (over 100M parameters), even though they achieve high performance. However, reproducing these methods would require considerable time and implementation effort. Therefore, due to time constraints, we don’t adapt these methods to our experimental setting.

To establish baselines, we adopt the existing retrieval results on the MSCOCO dataset reported in the respective papers of the compared methods and use them for comparison with our proposed approach. It should be noted that DynaCLIP is the result of transferring the DynaBERT (Hou et al., 2020) method to the CLIP framework, and the reported performance is sourced from MoPE-CLIP (Lin et al., 2024).

#### A.5 DETAILED TRAINING SETTINGS

Tab. 8 and Tab. 9 display the detailed training hyperparameters to get three scales of CLIP-Map models on YFCC15M dataset. Considering the smaller model capacity and weaker representational ability of CLIP-Map<sub>small</sub> and CLIP-Map<sub>tiny</sub>, we employ a longer mapping stage. For the larger CLIP-Map<sub>base</sub> variant, a shorter mapping stage proves adequate. All models are trained on 32 Nvidia H800 and it takes around 15 hours to finish the mapping-retraining pipeline for these three variant models each.

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

Table 8: Hyperparameters used in the Mapping Stage for different CLIP-Map model sizes.

Hyperparameter	CLIP-Map <sub>tiny</sub>	CLIP-Map <sub>small</sub>	CLIP-Map <sub>base</sub>
Batchsize	4096	4096	4096
Epochs	5	5	0.28(1000steps)
Distillation	False	False	False
Optimizer	AdamW	AdamW	AdamW
Optimizer Momentum	$\beta_1 = 0.9, \beta_2 = 0.98$	$\beta_1 = 0.9, \beta_2 = 0.98$	$\beta_1 = 0.9, \beta_2 = 0.98$
Base Learning Rate	$8 \times 10^{-4}$	$8 \times 10^{-4}$	$2 \times 10^{-3}$
Weight Decay	0.2	0.2	0.2
Learning Rate Schedule	Cosine decay	Cosine decay	Cosine decay
Warmup (steps)	2000	2000	1000
Gradient Clipping Norm	5	5	5
Reciprocal of Temperature	50	50	50
Image Resolution	$224 \times 224$	$224 \times 224$	$224 \times 224$
Image Augmentation	RandomResizedCrop	RandomResizedCrop	RandomResizedCrop
Tokenizer	Byte Pair Encoding	Byte Pair Encoding	Byte Pair Encoding
Vocabulary Size	49,408	49,408	49,408
Max Sequence Length	77	77	77

Table 9: Hyperparameters used in the Retraining Stage for different CLIP-Map model sizes.

Hyperparameter	CLIP-Map <sub>tiny</sub>	CLIP-Map <sub>small</sub>	CLIP-Map <sub>base</sub>
Batchsize	4096	4096	8192
Epochs	25	25	25
Optimizer	AdamW	AdamW	AdamW
Optimizer Momentum	$\beta_1 = 0.9, \beta_2 = 0.98$	$\beta_1 = 0.9, \beta_2 = 0.98$	$\beta_1 = 0.9, \beta_2 = 0.98$
Base Learning Rate	$1 \times 10^{-4}$	$1 \times 10^{-4}$	$3 \times 10^{-4}$
Weight Decay	0.2	0.2	$1 \times 10^{-4}$
Learning Rate Schedule	Cosine decay	Cosine decay	Cosine decay
Warmup (steps)	2000	2000	2000
Gradient Clipping Norm	5	5	5
Reciprocal of Temperature	50	50	50
Image Resolution	$224 \times 224$	$224 \times 224$	$224 \times 224$
Image Augmentation	RandomResizedCrop	RandomResizedCrop	RandomResizedCrop
Tokenizer	Byte Pair Encoding	Byte Pair Encoding	Byte Pair Encoding
Vocabulary Size	49,408	49,408	49,408
Max Sequence Length	77	77	77

Table 10: Performance comparison under different  $\lambda$  values.

$\lambda$	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
MSCOCO TR@1	33.1	37.2	40.1	42.2	43.8	45.6	48.1	49.4	50.2	<b>51.1</b>
MSCOCO IR@1	23.4	26.2	27.7	29.0	30.6	31.6	33.2	34.0	<b>34.7</b>	<b>34.7</b>

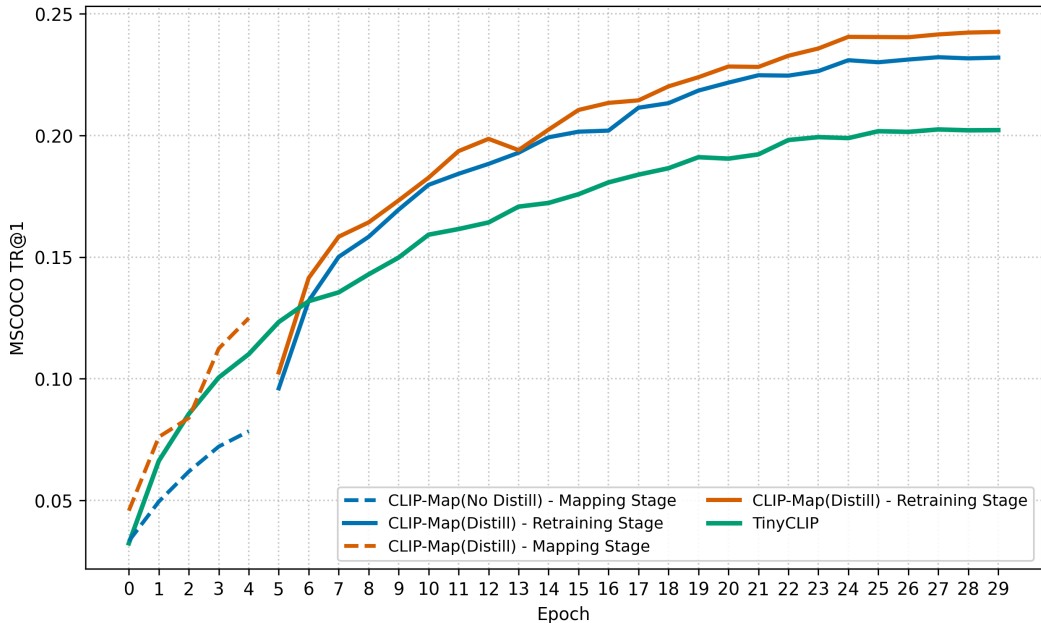


Figure 4: Evolution of TR@1 on the MSCOCO test set using TinyCLIP and CLIP-Map(w and w/o distillation during mapping stage).

## A.6 TRAINING AND INFERENCE SPEED-UP

Table 11: Wall-clock training time and epochs for different model variants.

Model	Wall-Clock Time (Multi-stages)	Epochs
TinyCLIP <sub>base</sub>	14h43m	25ep (25ep)
<b>CLIP-Map<sub>base</sub></b>	<b>14h30m</b> (6m + 14h23m)	25ep (1000steps + 25ep)
TinyCLIP <sub>small</sub>	32h36m (14h43m + 17h53m)	50ep (25ep + 25ep)
<b>CLIP-Map<sub>small</sub></b>	<b>22h10m</b> (4h50m + 17h19m)	25ep (5ep + 25ep)
TinyCLIP <sub>tiny</sub>	49h33m (14h43m + 17h53m + 16h57m)	75ep (25ep + 25ep + 25ep)
<b>CLIP-Map<sub>tiny</sub></b>	<b>21h50m</b> (4h37m + 17h11m)	25ep (5ep + 25ep)

We record the wall-clock time required for training different model variants, with the results presented in Tab. 11. Meanwhile, we also visualize the evolution of text-to-image retrieval performance (TR@1) on the MSCOCO test set during training, as shown in Fig. 4. Specifically, we compare our CLIP-Map method against TinyCLIP under the same compression setting (compressing OpenCLIP-ViT-B/16 to the *small* variant). Moreover, we investigate whether incorporating distillation during the mapping stage can further enhance performance. The results demonstrate that our method achieves faster convergence and consistently outperforms the baseline throughout training. It is worth noting that although incorporating distillation during the mapping stage can lead to a slight performance gain, we observe a drop in accuracy at the beginning of the retraining stage. We argue that such an instability may hinder subsequent optimization. Therefore, in our experiments, we avoid using distillation during the mapping stage in all subsequent experiments.



Table 12: Comparison of computational cost(GFLOPs/pair) and throughput(Pairs/s) across different model variants.

Model	GFLOPs (img + text)	Throughput (Pairs/s)
ViT-B/16	17.56 + 2.98	1041 pairs/s
CLIP-Map <sub>base</sub> / TinyCLIP <sub>base</sub>	7.99 + 1.49	1654 pairs/s ( $\times 1.6$ )
CLIP-Map <sub>small</sub> / TinyCLIP <sub>small</sub>	1.79 + 0.19	2466 pairs/s ( $\times 2.4$ )
CLIP-Map <sub>tiny</sub> / TinyCLIP <sub>tiny</sub>	0.21 + 0.03	3847 pairs/s ( $\times 3.8$ )

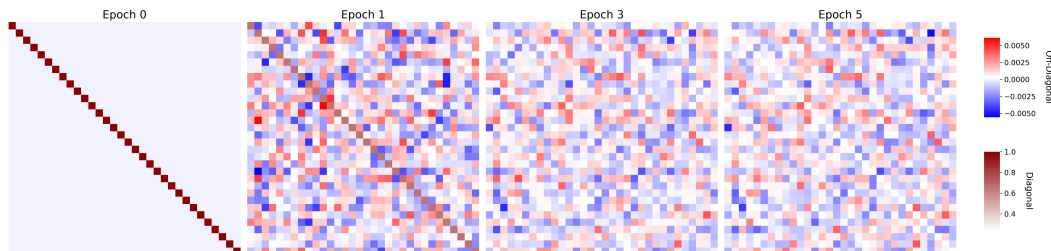


Figure 5: Changes of a mapping matrix in CLIP-Map<sub>small</sub> mapping stage. Due to the significant scale difference between diagonal and off-diagonal weights, we adopt two separate color scales to enhance the visualization quality.

The enhancement in the inference efficiency of our model is measured using the throughput of image–text pairs and the GFLOPs cost per image-text pair. All metrics are tested on single H800 GPU, at a batch size of 32 with an image resolution of 224×224, shown in Tab. 12. Because our compressed model shares the same parameter count and architecture as its TinyCLIP counterpart of equivalent size, the FLOPs and throughput measurements are identical for both.

#### A.7 VISUALIZATION

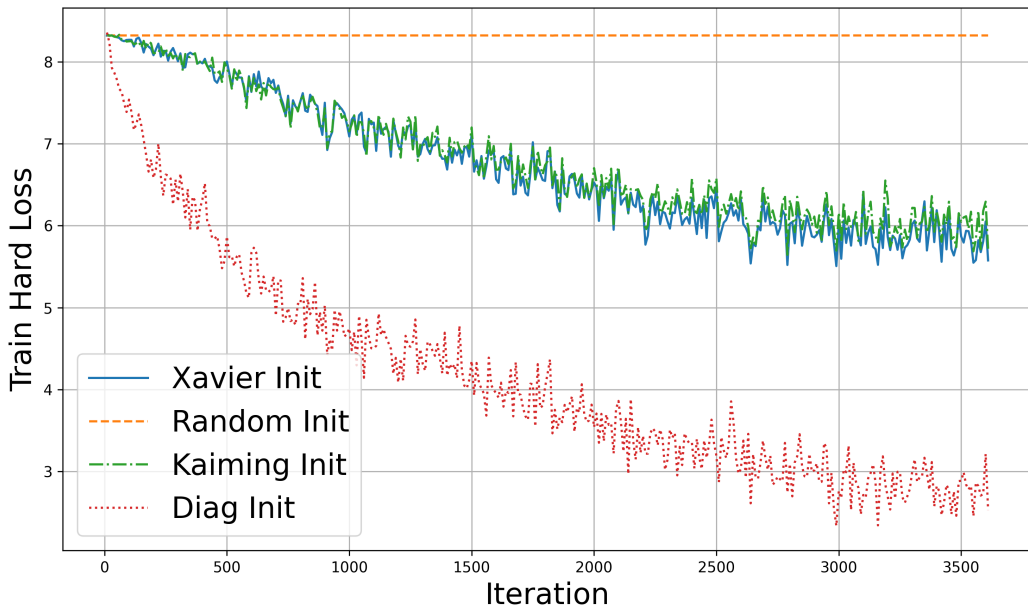


Figure 6: Loss curves under different initialization methods.

918 To better understand the mapping stage and the effect of different initialization methods, we visualize  
919 the changes in the mapping matrices across different epochs during the mapping stage of CLIP-  
920 Map<sub>small</sub> training in Fig. 5 and loss curves under different initialization methods during mapping  
921 stage in Fig. 6. For clarity, we only present a subset of the parameter matrix.  
922

#### 923 A.8 ABLATION ON TRAINING LOSS

924  
925 To efficiently identify the optimal distillation coefficient  $\lambda$ , we conduct a series of experiments on  
926 the small-scale CC3M dataset. In this setting, we compress the OpenCLIP-ViT-B/16 model into  
927 CLIP-Map<sub>small</sub>, and evaluate performance on MSCOCO test set under various  $\lambda$  values. The results  
928 are reported in Tab. 10. As the value of  $\lambda$  increases, the distillation loss becomes more dominant  
929 during training, leading to consistent performance improvements. Notably, when  $\lambda = 1.0$ , i.e., the  
930 total loss is fully composed of the distillation objective ( $\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{distill}}$ ), the model achieves its best  
931 performance. Therefore, we adopt  $\lambda = 1.0$  as the default setting in all subsequent experiments.

#### 932 A.9 LIMITATIONS AND FUTURE WORKS

933  
934 In this paper, we propose a mapping-based compression method termed CLIP-Map, using learnable  
935 matrices to mapping the original model into a smaller one. Our CLIP-Map outperforms baseline  
936 methods under equivalent compression settings, demonstrating both effectiveness and efficiency.  
937 However, our study is constrained by limited computational resources, which prevented us from  
938 evaluating the method on larger-scale datasets.

939 In future work, we plan to extend our approach to larger training datasets such as LAION-2B (Schuh-  
940 mann et al., 2022) to further validate its effectiveness and generalizability.  
941

#### 942 A.10 BROADER IMPACTS

943  
944 This work focuses on model compression techniques for CLIP, with the aim of enabling efficient  
945 deployment on resource-constrained devices. Since CLIP primarily performs vision-language under-  
946 standing rather than generation, it is less likely to introduce negative societal impacts. To the best of  
947 our knowledge, our paper does not pose any negative societal impact.  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971